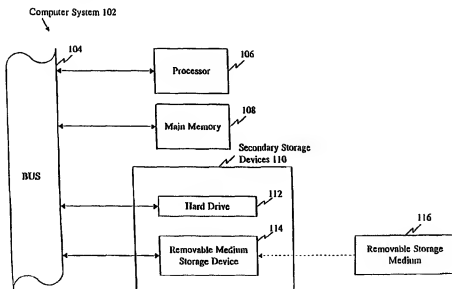




## INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification <sup>6</sup> : <b>C12N 15/31, C07K 14/315, 16/12, C12Q 1/68</b>		(11) International Publication Number: <b>WO 98/18931</b>	
<b>A2</b>		(43) International Publication Date: 7 May 1998 (07.05.98)	
(21) International Application Number: PCT/US97/19588		(74) Agents: BROOKES, A., Anders et al.; Human Genome Sciences, Inc., 9410 Key West Avenue, Rockville, MD 20850 (US).	
(22) International Filing Date: 30 October 1997 (30.10.97)			
(30) Priority Data: 60/029,960      31 October 1996 (31.10.96)      US	(81) Designated States: AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GE, GH, HU, ID, IL, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZW, ARIPO patent (GH, KE, LS, MW, SD, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, ML, MR, NE, SN, TD, TG).		
(71) Applicant (for all designated States except US): HUMAN GENOME SCIENCES, INC. [US/US]; 9410 Key West Avenue, Rockville, MD 20850 (US).			
(72) Inventors; and (75) Inventors/Applicants (for US only): KUNSCH, Charles, A. [US/US]; 2398B Dunwoody Crossing, Atlanta, GA 30338 (US). CHIOI, Gil, H. [KR/US]; 11429 Potomac Oaks Drive, Rockville, MD 20850 (US). DILLON, Patrick, J. [US/US]; 1055 Snipe Court, Carlsbad, CA 92009 (US). ROSEN, Craig, A. [US/US]; 22400 Rolling Hill Road, Laytonville, MD 20882 (US). BARASHI, Steven, C. [US/US]; 582 College Parkway #303, Rockville, MD 20850 (US). FANNON, Michael [US/US]; 13501 Rippling Brook Drive, Silver Spring, MD 20850 (US). DOUGHERTY, Brian, A. [US/US]; 708 Meadow Field Court, Mount Airy, MD 21771 (US).			
		<p><b>Published</b></p> <p>Without international search report and to be republished upon receipt of that report.</p>	

(54) Title: *STREPTOCOCCUS PNEUMONIAE* POLYNUCLEOTIDES AND SEQUENCES

## (57) Abstract

The present invention provides polynucleotide sequences of the genome of *Streptococcus pneumoniae*, polypeptide sequences encoded by the polynucleotide sequences, corresponding polynucleotides and polypeptides, vectors and hosts comprising the polynucleotides, and assays and other uses thereof. The present invention further provides polynucleotide and polypeptide sequence information stored on computer readable media, and computer-based systems and methods which facilitate its use.

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece			TR	Turkey
BG	Bulgaria	HU	Hungary	ML	Mali	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MN	Mongolia	UA	Ukraine
BR	Brazil	IL	Israel	MR	Mauritania	UG	Uganda
BY	Belarus	IS	Iceland	MW	Malawi	US	United States of America
CA	Canada	IT	Italy	MX	Mexico	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NE	Niger	VN	Viet Nam
CG	Congo	KE	Kenya	NL	Netherlands	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NO	Norway	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	NZ	New Zealand		
CM	Cameroon			PL	Poland		
CN	China	KR	Republic of Korea	PT	Portugal		
CU	Cuba	KZ	Kazakhstan	RO	Romania		
CZ	Czech Republic	LC	Saint Lucia	RU	Russian Federation		
DE	Germany	LI	Liechtenstein	SD	Sudan		
DK	Denmark	LK	Sri Lanka	SE	Sweden		
EE	Estonia	LR	Liberia	SG	Singapore		

## ***Streptococcus pneumoniae* Polynucleotides and Sequences**

### **FIELD OF THE INVENTION**

5           The present invention relates to the field of molecular biology. In particular, it relates to, among other things, nucleotide sequences of *Streptococcus pneumoniae*, contigs, ORFs, fragments, probes, primers and related polynucleotides thereof, peptides and polypeptides encoded by the sequences, and uses of the polynucleotides and sequences thereof, such as in fermentation, 10 polypeptide production, assays and pharmaceutical development, among others.

### **BACKGROUND OF THE INVENTION**

*Streptococcus pneumoniae* has been one of the most extensively studied 15 microorganisms since its first isolation in 1881. It was the object of many investigations that led to important scientific discoveries. In 1928, Griffith observed that when heat-killed encapsulated pneumococci and live strains constitutively lacking any capsule were concomitantly injected into mice, the nonencapsulated could be converted into encapsulated pneumococci with the same 20 capsular type as the heat-killed strain. Years later, the nature of this "transforming principle," or carrier of genetic information, was shown to be DNA. (Avery, O.T., et al., *J. Exp. Med.*, 79:137-157 (1944)).

          In spite of the vast number of publications on *S. pneumoniae* many questions about its virulence are still unanswered, and this pathogen remains a 25 major causative agent of serious human disease, especially community-acquired pneumonia. (Johnston, R.B., et al., *Rev. Infect. Dis.* 13(Suppl. 6):S509-517 (1991)). In addition, in developing countries, the pneumococcus is responsible for the death of a large number of children under the age of 5 years from pneumococcal pneumonia. The incidence of pneumococcal disease is highest in infants under 2 30 years of age and in people over 60 years of age. Pneumococci are the second most frequent cause (after *Haemophilus influenzae* type b) of bacterial meningitis and otitis media in children. With the recent introduction of conjugate vaccines for *H. influenzae* type b, pneumococcal meningitis is likely to become increasingly prominent. *S. pneumoniae* is the most important etiologic agent of community-

acquired pneumonia in adults and is the second most common cause of bacterial meningitis behind *Neisseria meningitidis*.

The antibiotic generally prescribed to treat *S. pneumoniae* is benzylpenicillin, although resistance to this and to other antibiotics is found occasionally. Pneumococcal resistance to penicillin results from mutations in its penicillin-binding proteins. In uncomplicated pneumococcal pneumonia caused by a sensitive strain, treatment with penicillin is usually successful unless started too late. Erythromycin or clindamycin can be used to treat pneumonia in patients hypersensitive to penicillin, but resistant strains to these drugs exist. Broad spectrum antibiotics (e.g., the tetracyclines) may also be effective, although tetracycline-resistant strains are not rare. In spite of the availability of antibiotics, the mortality of pneumococcal bacteremia in the last four decades has remained stable between 25 and 29%. (Gillespie, S.H., et al., *J. Med. Microbiol.* 28:237-248 (1989).

*S. pneumoniae* is carried in the upper respiratory tract by many healthy individuals. It has been suggested that attachment of pneumococci is mediated by a disaccharide receptor on fibronectin, present on human pharyngeal epithelial cells. (Anderson, B.J., et al., *J. Immunol.* 142:2464-2468 (1989). The mechanisms by which pneumococci translocate from the nasopharynx to the lung, thereby causing pneumonia, or migrate to the blood, giving rise to bacteremia or septicemia, are poorly understood. (Johnston, R.B., et al., *Rev. Infect. Dis.* 13(Suppl. 6):S509-517 (1991).

Various proteins have been suggested to be involved in the pathogenicity of *S. pneumoniae*, however, only a few of them have actually been confirmed as virulence factors. Pneumococci produce an IgA1 protease that might interfere with host defense at mucosal surfaces. (Kornfield, S.J., et al., *Rev. Inf. Dis.* 3:521-534 (1981). *S. pneumoniae* also produces neuraminidase, an enzyme that may facilitate attachment to epithelial cells by cleaving sialic acid from the host glycolipids and gangliosides. Partially purified neuraminidase was observed to induce meningitis-like symptoms in mice; however, the reliability of this finding has been questioned because the neuraminidase preparations used were probably contaminated with cell wall products. Other pneumococcal proteins besides neuraminidase are involved in the adhesion of pneumococci to epithelial and endothelial cells. These pneumococcal proteins have as yet not been identified. Recently, Cundell et al., reported that peptide permeases can modulate



pneumococcal adherence to epithelial and endothelial cells. It was, however, unclear whether these permeases function directly as adhesions or whether they enhance adherence by modulating the expression of pneumococcal adhesions. (DeVelasco, E.A., *et al.*, *Micro. Rev.* 59:591-603 (1995). A better understanding of the virulence factors determining its pathogenicity will need to be developed to cope with the devastating effects of pneumococcal disease in humans.

Ironically, despite the prominent role of *S. pneumoniae* in the discovery of DNA, little is known about the molecular genetics of the organism. The *S. pneumoniae* genome consists of one circular, covalently closed, double-stranded DNA and a collection of so-called variable accessory elements, such as prophages, plasmids, transposons and the like. Most physical characteristics and almost all of the genes of *S. pneumoniae* are unknown. Among the few that have been identified, most have not been physically mapped or characterized in detail. Only a few genes of this organism have been sequenced. (See, for instance current versions of GENBANK and other nucleic acid databases, and references that relate to the genome of *S. pneumoniae* such as those set out elsewhere herein.)

It is clear that the etiology of diseases mediated or exacerbated by *S. pneumoniae*, infection involves the programmed expression of *S. pneumoniae* genes, and that characterizing the genes and their patterns of expression would add dramatically to our understanding of the organism and its host interactions. Knowledge of *S. pneumoniae* genes and genomic organization would improve our understanding of disease etiology and lead to improved and new ways of preventing, ameliorating, arresting and reversing diseases. Moreover, characterized genes and genomic fragments of *S. pneumoniae* would provide reagents for, among other things, detecting, characterizing and controlling *S. pneumoniae* infections. There is a need to characterize the genome of *S. pneumoniae* and for polynucleotides of this organism.

### SUMMARY OF THE INVENTION

5 The present invention is based on the sequencing of fragments of the *Streptococcus pneumoniae* genome. The primary nucleotide sequences which were generated are provided in SEQ ID NOS:1-391.

10 The present invention provides the nucleotide sequence of several hundred contigs of the *Streptococcus pneumoniae* genome, which are listed in tables below and set out in the Sequence Listing submitted herewith, and representative fragments thereof, in a form which can be readily used, analyzed, and interpreted by a skilled artisan. In one embodiment, the present invention is provided as contiguous strings of primary sequence information corresponding to the nucleotide sequences depicted in SEQ ID NOS:1-391.

15 The present invention further provides nucleotide sequences which are at least 95% identical to the nucleotide sequences of SEQ ID NOS:1-391.

20 The nucleotide sequence of SEQ ID NOS:1-391, a representative fragment thereof, or a nucleotide sequence which is at least 95% identical to the nucleotide sequence of SEQ ID NOS:1-391 may be provided in a variety of mediums to facilitate its use. In one application of this embodiment, the sequences of the present invention are recorded on computer readable media. Such media includes, but is not limited to: magnetic storage media, such as floppy discs, hard disc storage medium, and magnetic tape; optical storage media such as CD-ROM; electrical storage media such as RAM and ROM; and hybrids of these categories such as magnetic/optical storage media.

25 The present invention further provides systems, particularly computer-based systems which contain the sequence information herein described stored in a data storage means. Such systems are designed to identify commercially important fragments of the *Streptococcus pneumoniae* genome.

30 Another embodiment of the present invention is directed to fragments of the *Streptococcus pneumoniae* genome having particular structural or functional attributes. Such fragments of the *Streptococcus pneumoniae* genome of the present invention include, but are not limited to, fragments which encode peptides, hereinafter referred to as open reading frames or ORFs, fragments which modulate the expression of an operably linked ORF, hereinafter referred to as expression modulating fragments or EMFs, and fragments which can be used to diagnose the

35

presence of *Streptococcus pneumoniae* in a sample, hereinafter referred to as diagnostic fragments or DFs.

Each of the ORFs in fragments of the *Streptococcus pneumoniae* genome disclosed in Tables 1-3, and the EMFs found 5' to the ORFs, can be used in numerous ways as polynucleotide reagents. For instance, the sequences can be used as diagnostic probes or amplification primers for detecting or determining the presence of a specific microbe in a sample, to selectively control gene expression in a host and in the production of polypeptides, such as polypeptides encoded by ORFs of the present invention, particular those polypeptides that have a pharmacological activity.

The present invention further includes recombinant constructs comprising one or more fragments of the *Streptococcus pneumoniae* genome of the present invention. The recombinant constructs of the present invention comprise vectors, such as a plasmid or viral vector, into which a fragment of the *Streptococcus pneumoniae* has been inserted.

The present invention further provides host cells containing any of the isolated fragments of the *Streptococcus pneumoniae* genome of the present invention. The host cells can be a higher eukaryotic host cell, such as a mammalian cell, a lower eukaryotic cell, such as a yeast cell, or a procaryotic cell such as a bacterial cell.

The present invention is further directed to isolated polypeptides and proteins encoded by ORFs of the present invention. A variety of methods, well known to those of skill in the art, routinely may be utilized to obtain any of the polypeptides and proteins of the present invention. For instance, polypeptides and proteins of the present invention having relatively short, simple amino acid sequences readily can be synthesized using commercially available automated peptide synthesizers. Polypeptides and proteins of the present invention also may be purified from bacterial cells which naturally produce the protein. Yet another alternative is to purify polypeptide and proteins of the present invention from cells which have been altered to express them.

The invention further provides methods of obtaining homologs of the fragments of the *Streptococcus pneumoniae* genome of the present invention and homologs of the proteins encoded by the ORFs of the present invention. Specifically, by using the nucleotide and amino acid sequences disclosed herein as

a probe or as primers, and techniques such as PCR cloning and colony/plaque hybridization, one skilled in the art can obtain homologs.

The invention further provides antibodies which selectively bind polypeptides and proteins of the present invention. Such antibodies include both  
5 monoclonal and polyclonal antibodies.

The invention further provides hybridomas which produce the above-described antibodies. A hybridoma is an immortalized cell line which is capable of secreting a specific monoclonal antibody.

The present invention further provides methods of identifying test samples  
10 derived from cells which express one of the ORFs of the present invention, or a homolog thereof. Such methods comprise incubating a test sample with one or more of the antibodies of the present invention, or one or more of the DFs of the present invention, under conditions which allow a skilled artisan to determine if the sample contains the ORF or product produced therefrom.

15 In another embodiment of the present invention, kits are provided which contain the necessary reagents to carry out the above-described assays.

Specifically, the invention provides a compartmentalized kit to receive, in close confinement, one or more containers which comprises: (a) a first container comprising one of the antibodies, or one of the DFs of the present invention; and  
20 (b) one or more other containers comprising one or more of the following: wash reagents, reagents capable of detecting presence of bound antibodies or hybridized DFs.

Using the isolated proteins of the present invention, the present invention further provides methods of obtaining and identifying agents capable of binding to  
25 a polypeptide or protein encoded by one of the ORFs of the present invention. Specifically, such agents include, as further described below, antibodies, peptides, carbohydrates, pharmaceutical agents and the like. Such methods comprise steps of: (a) contacting an agent with an isolated protein encoded by one of the ORFs of the present invention; and (b) determining whether the agent binds to said protein.

30 The present genomic sequences of *Streptococcus pneumoniae* will be of great value to all laboratories working with this organism and for a variety of commercial purposes. Many fragments of the *Streptococcus pneumoniae* genome will be immediately identified by similarity searches against GenBank or protein databases and will be of immediate value to *Streptococcus pneumoniae* researchers

and for immediate commercial value for the production of proteins or to control gene expression.

The methodology and technology for elucidating extensive genomic sequences of bacterial and other genomes has and will greatly enhance the ability to analyze and understand chromosomal organization. In particular, sequenced contigs and genomes will provide the models for developing tools for the analysis of chromosome structure and function, including the ability to identify genes within large segments of genomic DNA, the structure, position, and spacing of regulatory elements, the identification of genes with potential industrial applications, and the ability to do comparative genomic and molecular phylogeny.

### DESCRIPTION OF THE FIGURES

FIGURE 1 is a block diagram of a computer system (102) that can be used to implement computer-based systems of present invention.

FIGURE 2 is a schematic diagram depicting the data flow and computer programs used to collect, assemble, edit and annotate the contigs of the *Streptococcus pneumoniae* genome of the present invention. Both Macintosh and Unix platforms are used to handle the AB 373 and 377 sequence data files, largely as described in Kerlavage *et al.*, *Proceedings of the Twenty-Sixth Annual Hawaii International Conference on System Sciences*, 585, IEEE Computer Society Press, Washington D.C. (1993). Factura (AB) is a Macintosh program designed for automatic vector sequence removal and end-trimming of sequence files. The program Loadis runs on a Macintosh platform and parses the feature data extracted from the sequence files by Factura to the Unix based *Streptococcus pneumoniae* relational database. Assembly of contigs (and whole genome sequences) is accomplished by retrieving a specific set of sequence files and their associated features using Extrseq, a Unix utility for retrieving sequences from an SQL database. The resulting sequence file is processed by seq\_filter to trim portions of the sequences with more than 2% ambiguous nucleotides. The sequence files were assembled using TIGR Assembler, an assembly engine designed at The Institute for Genomic Research (TIGR) for rapid and accurate assembly of thousands of sequence fragments. The collection of contigs generated by the assembly step is loaded into the database with the lassie program. Identification of open reading

frames (ORFs) is accomplished by processing contigs with zorf or GenMark. The ORFs are searched against *S. pneumoniae* sequences from GenBank and against all protein sequences using the BLASTN and BLASTP programs, described in Altschul *et al.*, *J. Mol. Biol.* 215: 403-410 (1990)). Results of the ORF  
5 determination and similarity searching steps were loaded into the database. As described below, some results of the determination and the searches are set out in Tables 1-3.

### DETAILED DESCRIPTION OF ILLUSTRATIVE EMBODIMENTS

10 The present invention is based on the sequencing of fragments of the *Streptococcus pneumoniae* genome and analysis of the sequences. The primary nucleotide sequences generated by sequencing the fragments are provided in SEQ ID NOS:1-391. (As used herein, the "primary sequence" refers to the nucleotide  
15 sequence represented by the IUPAC nomenclature system.)

In addition to the aforementioned *Streptococcus pneumoniae* polynucleotide and polynucleotide sequences, the present invention provides the nucleotide sequences of SEQ ID NOS:1-391, or representative fragments thereof, in a form which can be readily used, analyzed, and interpreted by a skilled artisan.

20 As used herein, a "representative fragment of the nucleotide sequence depicted in SEQ ID NOS:1-391" refers to any portion of the SEQ ID NOS:1-391 which is not presently represented within a publicly available database. Preferred representative fragments of the present invention are *Streptococcus pneumoniae* open reading frames ( ORFs ), expression modulating fragment ( EMFs ) and  
25 fragments which can be used to diagnose the presence of *Streptococcus pneumoniae* in sample ( DFs ). A non-limiting identification of preferred representative fragments is provided in Tables 1-3. As discussed in detail below, the information provided in SEQ ID NOS:1-391 and in Tables 1-3 together with routine cloning, synthesis, sequencing and assay methods will enable those skilled  
30 in the art to clone and sequence all "representative fragments" of interest, including open reading frames encoding a large variety of *Streptococcus pneumoniae* proteins.

While the presently disclosed sequences of SEQ ID NOS:1-391 are highly accurate, sequencing techniques are not perfect and, in relatively rare instances,  
35 further investigation of a fragment or sequence of the invention may reveal a

nucleotide sequence error present in a nucleotide sequence disclosed in SEQ ID NOS:1-391. However, once the present invention is made available (*i.e.*, once the information in SEQ ID NOS:1-391 and Tables 1-3 has been made available), resolving a rare sequencing error in SEQ ID NOS:1-391 will be well within the skill of the art. The present disclosure makes available sufficient sequence information to allow any of the described contigs or portions thereof to be obtained readily by straightforward application of routine techniques. Further sequencing of such polynucleotide may proceed in like manner using manual and automated sequencing methods which are employed ubiquitous in the art. Nucleotide sequence editing software is publicly available. For example, Applied Biosystem's (AB) AutoAssembler can be used as an aid during visual inspection of nucleotide sequences. By employing such routine techniques potential errors readily may be identified and the correct sequence then may be ascertained by targeting further sequencing effort, also of a routine nature, to the region containing the potential error.

Even if all of the very rare sequencing errors in SEQ ID NOS:1-391 were corrected, the resulting nucleotide sequences would still be at least 95% identical, nearly all would be at least 99% identical, and the great majority would be at least 99.9% identical to the nucleotide sequences of SEQ ID NOS:1-391.

As discussed elsewhere herein, polynucleotides of the present invention readily may be obtained by routine application of well known and standard procedures for cloning and sequencing DNA. Detailed methods for obtaining libraries and for sequencing are provided below, for instance. A wide variety of *Streptococcus pneumoniae* strains that can be used to prepare *S. pneumoniae* genomic DNA for cloning and for obtaining polynucleotides of the present invention are available to the public from recognized depository institutions, such as the American Type Culture Collection (ATCC). While the present invention is enabled by the sequences and other information herein disclosed, the *S. pneumoniae* strain that provided the DNA of the present Sequence Listing, Strain 78/7 14.8.91, has been deposited in the ATCC, as a convenience to those of skill in the art. As a further convenience, a library of *S. pneumoniae* genomic DNA, derived from the same strain, also has been deposited in the ATCC. The *S. pneumoniae* strain was deposited on October 10, 1996, and was given Deposit No. 55840, and the cDNA library was deposited on October 11, 1996 and was given Deposit No. 97755. The genomic fragments in the library are 15 to 20 kb

fragments generated by partial Sau3A1 digestion and they are inserted into the BamHI site in the well-known lambda-derived vector lambda DASH II (Stratagene, La Jolla, CA). The provision of the deposits is not a waiver of any rights of the inventors or their assignees in the present subject matter.

5 The nucleotide sequences of the genomes from different strains of *Streptococcus pneumoniae* differ somewhat. However, the nucleotide sequences of the genomes of all *Streptococcus pneumoniae* strains will be at least 95% identical, in corresponding part, to the nucleotide sequences provided in SEQ ID NOS:1-391. Nearly all will be at least 99% identical and the great majority will be  
10 99.9% identical.

Thus, the present invention further provides nucleotide sequences which are at least 95%, preferably 99% and most preferably 99.9% identical to the nucleotide sequences of SEQ ID NOS:1-391, in a form which can be readily used, analyzed and interpreted by the skilled artisan.

15 Methods for determining whether a nucleotide sequence is at least 95%, at least 99% or at least 99.9% identical to the nucleotide sequences of SEQ ID NOS:1-391 are routine and readily available to the skilled artisan. For example, the well known fasta algorithm described in Pearson and Lipman, *Proc. Natl. Acad. Sci. USA* 85: 2444 (1988) can be used to generate the percent identity of nucleotide  
20 sequences. The BLASTN program also can be used to generate an identity score of polynucleotides compared to one another.

## COMPUTER RELATED EMBODIMENTS

The nucleotide sequences provided in SEQ ID NOS:1-391, a representative  
25 fragment thereof, or a nucleotide sequence at least 95%, preferably at least 99% and most preferably at least 99.9% identical to a polynucleotide sequence of SEQ ID NOS:1-391 may be "provided" in a variety of mediums to facilitate use thereof. As used herein, provided refers to a manufacture, other than an isolated nucleic acid molecule, which contains a nucleotide sequence of the present invention; i.e.,  
30 a nucleotide sequence provided in SEQ ID NOS:1-391, a representative fragment thereof, or a nucleotide sequence at least 95%, preferably at least 99% and most preferably at least 99.9% identical to a polynucleotide of SEQ ID NOS:1-391. Such a manufacture provides a large portion of the *Streptococcus pneumoniae* genome and parts thereof (e.g., a *Streptococcus pneumoniae* open reading frame  
35 (ORF)) in a form which allows a skilled artisan to examine the manufacture using



means not directly applicable to examining the *Streptococcus pneumoniae* genome or a subset thereof as it exists in nature or in purified form.

In one application of this embodiment, a nucleotide sequence of the present invention can be recorded on computer readable media. As used herein, "computer readable media" refers to any medium which can be read and accessed directly by a computer. Such media include, but are not limited to: magnetic storage media, such as floppy discs, hard disc storage medium, and magnetic tape; optical storage media such as CD-ROM; electrical storage media such as RAM and ROM; and hybrids of these categories, such as magnetic/optical storage media. A skilled artisan can readily appreciate how any of the presently known computer readable mediums can be used to create a manufacture comprising computer readable medium having recorded thereon a nucleotide sequence of the present invention. Likewise, it will be clear to those of skill how additional computer readable media that may be developed also can be used to create analogous manufactures having recorded thereon a nucleotide sequence of the present invention.

As used herein, "recorded" refers to a process for storing information on computer readable medium. A skilled artisan can readily adopt any of the presently known methods for recording information on computer readable medium to generate manufactures comprising the nucleotide sequence information of the present invention. A variety of data storage structures are available to a skilled artisan for creating a computer readable medium having recorded thereon a nucleotide sequence of the present invention. The choice of the data storage structure will generally be based on the means chosen to access the stored information. In addition, a variety of data processor programs and formats can be used to store the nucleotide sequence information of the present invention on computer readable medium. The sequence information can be represented in a word processing text file, formatted in commercially-available software such as WordPerfect and Microsoft Word, or represented in the form of an ASCII file, stored in a database application, such as DB2, Sybase, Oracle, or the like. A skilled artisan can readily adapt any number of data-processor structuring formats (e.g., text file or database) in order to obtain computer readable medium having recorded thereon the nucleotide sequence information of the present invention.

Computer software is publicly available which allows a skilled artisan to access sequence information provided in a computer readable medium. Thus, by providing in computer readable form the nucleotide sequences of SEQ ID NOS:1-

391, a representative fragment thereof, or a nucleotide sequence at least 95%, preferably at least 99% and most preferably at least 99.9% identical to a sequence of SEQ ID NOS:1-391 the present invention enables the skilled artisan routinely to access the provided sequence information for a wide variety of purposes.

5       The examples which follow demonstrate how software which implements the BLAST (Altschul *et al.*, *J. Mol. Biol.* 215:403-410 (1990)) and BLAZE (Brutlag *et al.*, *Comp. Chem.* 17:203-207 (1993)) search algorithms on a Sybase system was used to identify open reading frames (ORFs) within the *Streptococcus pneumoniae* genome which contain homology to ORFs or proteins from both  
10 *Streptococcus pneumoniae* and from other organisms. Among the ORFs discussed herein are protein encoding fragments of the *Streptococcus pneumoniae* genome useful in producing commercially important proteins, such as enzymes used in fermentation reactions and in the production of commercially useful metabolites.

15       The present invention further provides systems, particularly computer-based systems, which contain the sequence information described herein. Such systems are designed to identify, among other things, commercially important fragments of the *Streptococcus pneumoniae* genome.

20       As used herein, "a computer-based system" refers to the hardware means, software means, and data storage means used to analyze the nucleotide sequence information of the present invention. The minimum hardware means of the computer-based systems of the present invention comprises a central processing unit (CPU), input means, output means, and data storage means. A skilled artisan can readily appreciate that any one of the currently available computer-based systems are suitable for use in the present invention.

25       As stated above, the computer-based systems of the present invention comprise a data storage means having stored therein a nucleotide sequence of the present invention and the necessary hardware means and software means for supporting and implementing a search means.

30       As used herein, "data storage means" refers to memory which can store nucleotide sequence information of the present invention, or a memory access means which can access manufactures having recorded thereon the nucleotide sequence information of the present invention.

35       As used herein, "search means" refers to one or more programs which are implemented on the computer-based system to compare a target sequence or target structural motif with the sequence information stored within the data storage

means. Search means are used to identify fragments or regions of the present genomic sequences which match a particular target sequence or target motif. A variety of known algorithms are disclosed publicly and a variety of commercially available software for conducting search means are and can be used in the computer-based systems of the present invention. Examples of such software includes, but is not limited to, MacPattern (EMBL), BLASTN and BLASTX (NCBIA). A skilled artisan can readily recognize that any one of the available algorithms or implementing software packages for conducting homology searches can be adapted for use in the present computer-based systems.

As used herein, a "target sequence" can be any DNA or amino acid sequence of six or more nucleotides or two or more amino acids. A skilled artisan can readily recognize that the longer a target sequence is, the less likely a target sequence will be present as a random occurrence in the database. The most preferred sequence length of a target sequence is from about 10 to 100 amino acids or from about 30 to 300 nucleotide residues. However, it is well recognized that searches for commercially important fragments, such as sequence fragments involved in gene expression and protein processing, may be of shorter length.

As used herein, "a target structural motif," or "target motif," refers to any rationally selected sequence or combination of sequences in which the sequence(s) are chosen based on a three-dimensional configuration which is formed upon the folding of the target motif. There are a variety of target motifs known in the art. Protein target motifs include, but are not limited to, enzymic active sites and signal sequences. Nucleic acid target motifs include, but are not limited to, promoter sequences, hairpin structures and inducible expression elements (protein binding sequences).

A variety of structural formats for the input and output means can be used to input and output the information in the computer-based systems of the present invention. A preferred format for an output means ranks fragments of the *Streptococcus pneumoniae* genomic sequences possessing varying degrees of homology to the target sequence or target motif. Such presentation provides a skilled artisan with a ranking of sequences which contain various amounts of the target sequence or target motif and identifies the degree of homology contained in the identified fragment.

A variety of comparing means can be used to compare a target sequence or target motif with the data storage means to identify sequence fragments of the

*Streptococcus pneumoniae* genome. In the present examples, implementing software which implement the BLAST and BLAZE algorithms, described in Altschul *et al.*, *J. Mol. Biol.* 215: 403-410 (1990), is used to identify open reading frames within the *Streptococcus pneumoniae* genome. A skilled artisan can readily recognize that any one of the publicly available homology search programs can be used as the search means for the computer-based systems of the present invention. Of course, suitable proprietary systems that may be known to those of skill also may be employed in this regard.

Figure 1 provides a block diagram of a computer system illustrative of embodiments of this aspect of present invention. The computer system 102 includes a processor 106 connected to a bus 104. Also connected to the bus 104 are a main memory 108 (preferably implemented as random access memory, RAM) and a variety of secondary storage devices 110, such as a hard drive 112 and a removable medium storage device 114. The removable medium storage device 114 may represent, for example, a floppy disk drive, a CD-ROM drive, a magnetic tape drive, *etc.* A removable storage medium 116 (such as a floppy disk, a compact disk, a magnetic tape, *etc.*) containing control logic and/or data recorded therein may be inserted into the removable medium storage device 114. The computer system 102 includes appropriate software for reading the control logic and/or the data from the removable medium storage device 114, once it is inserted into the removable medium storage device 114.

A nucleotide sequence of the present invention may be stored in a well known manner in the main memory 108, any of the secondary storage devices 110, and/or a removable storage medium 116. During execution, software for accessing and processing the genomic sequence (such as search tools, comparing tools, *etc.*) reside in main memory 108, in accordance with the requirements and operating parameters of the operating system, the hardware system and the software program or programs.

### BIOCHEMICAL EMBODIMENTS

Other embodiments of the present invention are directed to isolated fragments of the *Streptococcus pneumoniae* genome. The fragments of the  
5 *Streptococcus pneumoniae* genome of the present invention include, but are not limited to fragments which encode peptides and polypeptides, hereinafter open reading frames (ORFs), fragments which modulate the expression of an operably linked ORF, hereinafter expression modulating fragments (EMFs) and fragments which can be used to diagnose the presence of *Streptococcus pneumoniae* in a  
10 sample, hereinafter diagnostic fragments (DFs).

As used herein, an "isolated nucleic acid molecule" or an "isolated fragment of the *Streptococcus pneumoniae* genome" refers to a nucleic acid molecule possessing a specific nucleotide sequence which has been subjected to purification means to reduce, from the composition, the number of compounds which are  
15 normally associated with the composition. Particularly, the term refers to the nucleic acid molecules having the sequences set out in SEQ ID NOS:1-391, to representative fragments thereof as described above, to polynucleotides at least 95%, preferably at least 99% and especially preferably at least 99.9% identical in sequence thereto, also as set out above.

20 A variety of purification means can be used to generate the isolated fragments of the present invention. These include, but are not limited to methods which separate constituents of a solution based on charge, solubility, or size.

In one embodiment, *Streptococcus pneumoniae* DNA can be enzymatically sheared to produce fragments of 15-20 kb in length. These fragments can then be  
25 used to generate a *Streptococcus pneumoniae* library by inserting them into lambda clones as described in the Examples below. Primers flanking, for example, an ORF, such as those enumerated in Tables 1-3 can then be generated using nucleotide sequence information provided in SEQ ID NOS:1-391. Well known and routine techniques of PCR cloning then can be used to isolate the ORF from  
30 the lambda DNA library or *Streptococcus pneumoniae* genomic DNA. Thus, given the availability of SEQ ID NOS:1-391, the information in Tables 1, 2 and 3, and the information that may be obtained readily by analysis of the sequences of SEQ ID NOS:1-391 using methods set out above, those of skill will be enabled by the present disclosure to isolate any ORF-containing or other nucleic acid fragment of  
35 the present invention.

The isolated nucleic acid molecules of the present invention include, but are not limited to single stranded and double stranded DNA, and single stranded RNA.

As used herein, an "open reading frame," ORF, means a series of triplets coding for amino acids without any termination codons and is a sequence translatable into protein.

Tables 1, 2, and 3 list ORFs in the *Streptococcus pneumoniae* genomic contigs of the present invention that were identified as putative coding regions by the GeneMark software using organism-specific second-order Markov probability transition matrices. It will be appreciated that other criteria can be used, in accordance with well known analytical methods, such as those discussed herein, to generate more inclusive, more restrictive, or more selective lists.

Table 1 sets out ORFs in the *Streptococcus pneumoniae* contigs of the present invention that over a continuous region of at least 50 bases are 95% or more identical (by BLAST analysis) to a nucleotide sequence available through GenBank in October, 1997.

Table 2 sets out ORFs in the *Streptococcus pneumoniae* contigs of the present invention that are not in Table 1 and match, with a BLASTP probability score of 0.01 or less, a polypeptide sequence available through GenBank in October, 1997.

Table 3 sets out ORFs in the *Streptococcus pneumoniae* contigs of the present invention that do not match significantly, by BLASTP analysis, a polypeptide sequence available through GenBank in October, 1997.

In each table, the first and second columns identify the ORF by, respectively, contig number and ORF number within the contig; the third column indicates the first nucleotide of the ORF (actually the first nucleotide of the stop codon immediately preceding the ORF), counting from the 5' end of the contig strand; and the fourth column, "stop (nt)" indicates the last nucleotide of the stop codon defining the 3' end of the ORF.

In Tables 1 and 2, column five, lists the Reference for the closest matching sequence available through GenBank. These reference numbers are the databases entry numbers commonly used by those of skill in the art, who will be familiar with their denominators. Descriptions of the nomenclature are available from the National Center for Biotechnology Information. Column six in Tables 1 and 2 provides the gene name of the matching sequence; column seven provides the BLAST identity score and column eight the BLAST similarity score from the

comparison of the ORF and the homologous gene; and column nine indicates the length in nucleotides of the highest scoring segment pair identified by the BLAST identity analysis.

Each ORF described in the tables is defined by "start (nt)" (5') and "stop (nt)" (3') nucleotide position numbers. These position numbers refer to the boundaries of each ORF and provide orientation with respect to whether the forward or reverse strand is the coding strand and which reading frame the coding sequence is contained. The "start" position is the first nucleotide of the triplet encoding a stop codon just 5' to the ORF and the "stop" position is the last nucleotide of the triplet encoding the next in-frame stop codon (i.e., the stop codon at the 3' end of the ORF). Those of ordinary skill in the art appreciate that preferred fragments within each ORF described in the table include fragments of each ORF which include the entire sequence from the delineated "start" and "stop" positions excepting the first and last three nucleotides since these encode stop codons. Thus, polynucleotides set out as ORFs in the tables but lacking the three (3) 5' nucleotides and the three (3) 3' nucleotides are encompassed by the present invention. Those of skill also appreciate that particularly preferred are fragments within each ORF that are polynucleotide fragments comprising polypeptide coding sequence. As defined herein, "coding sequence" includes the fragment within an ORF beginning at the first in-frame ATG (triplet encoding methionine) and ending with the last nucleotide prior to the triplet encoding the 3' stop codon. Preferred are fragments comprising the entire coding sequence and fragments comprising the entire coding sequence, excepting the coding sequence for the N-terminal methionine. Those of skill appreciate that the N-terminal methionine is often removed during post-translational processing and that polynucleotides lacking the ATG can be used to facilitate production of N-terminal fusion proteins which may be beneficial in the production or use of genetically engineered proteins. Of course, due to the degeneracy of the genetic code many polynucleotides can encode a given polypeptide. Thus, the invention further includes polynucleotides comprising a nucleotide sequence encoding a polypeptide sequence itself encoded by the coding sequence within an ORF described in Tables 1-3 herein. Further, polynucleotides at least 95%, preferably at least 99% and especially preferably at least 99.9% identical in sequence to the foregoing polynucleotides, are contemplated by the present invention.

Polypeptides encoded by polynucleotides described above and elsewhere herein are also provided by the present invention as are polypeptide comprising an amino acid sequence at least about 95%, preferably at least 97% and even more preferably 99% identical to the amino acid sequence of a polypeptide encoded by an ORF shown in Tables 1-3. These polypeptides may or may not comprise an N-terminal methionine.

The concepts of percent identity and percent similarity of two polypeptide sequences is well understood in the art. For example, two polypeptides 10 amino acids in length which differ at three amino acid positions (*e.g.*, at positions 1, 3 and 5) are said to have a percent identity of 70%. However, the same two polypeptides would be deemed to have a percent similarity of 80% if, for example at position 5, the amino acids moieties, although not identical, were "similar" (*i.e.*, possessed similar biochemical characteristics). Many programs for analysis of nucleotide or amino acid sequence similarity, such as fasta and BLAST specifically list percent identity of a matching region as an output parameter. Thus, for instance, Tables 1 and 2 herein enumerate the percent identity of the highest scoring segment pair in each ORF and its listed relative. Further details concerning the algorithms and criteria used for homology searches are provided below and are described in the pertinent literature highlighted by the citations provided below.

It will be appreciated that other criteria can be used to generate more inclusive and more exclusive listings of the types set out in the tables. As those of skill will appreciate, narrow and broad searches both are useful. Thus, a skilled artisan can readily identify ORFs in contigs of the *Streptococcus pneumoniae* genome other than those listed in Tables 1-3, such as ORFs which are overlapping or encoded by the opposite strand of an identified ORF in addition to those ascertainable using the computer-based systems of the present invention.

As used herein, an "expression modulating fragment," EMF, means a series of nucleotide molecules which modulates the expression of an operably linked ORF or EMF.



As used herein, a sequence is said to "modulate the expression of an operably linked sequence" when the expression of the sequence is altered by the presence of the EMF. EMFs include, but are not limited to, promoters, and promoter modulating sequences (inducible elements). One class of EMFs are fragments which induce the expression or an operably linked ORF in response to a specific regulatory factor or physiological event.

EMF sequences can be identified within the contigs of the *Streptococcus pneumoniae* genome by their proximity to the ORFs provided in Tables 1-3. An intergenic segment, or a fragment of the intergenic segment, from about 10 to 200 nucleotides in length, taken from any one of the ORFs of Tables 1-3 will modulate the expression of an operably linked ORF in a fashion similar to that found with the naturally linked ORF sequence. As used herein, an "intergenic segment" refers to fragments of the *Streptococcus pneumoniae* genome which are between two ORF(s) herein described. EMFs also can be identified using known EMFs as a target sequence or target motif in the computer-based systems of the present invention. Further, the two methods can be combined and used together.

The presence and activity of an EMF can be confirmed using an EMF trap vector. An EMF trap vector contains a cloning site linked to a marker sequence. A marker sequence encodes an identifiable phenotype, such as antibiotic resistance or a complementing nutrition auxotrophic factor, which can be identified or assayed when the EMF trap vector is placed within an appropriate host under appropriate conditions. As described above, a EMF will modulate the expression of an operably linked marker sequence. A more detailed discussion of various marker sequences is provided below. A sequence which is suspected as being an EMF is cloned in all three reading frames in one or more restriction sites upstream from the marker sequence in the EMF trap vector. The vector is then transformed into an appropriate host using known procedures and the phenotype of the transformed host is examined under appropriate conditions. As described above, an EMF will modulate the expression of an operably linked marker sequence.

As used herein, a "diagnostic fragment," DF, means a series of nucleotide molecules which selectively hybridize to *Streptococcus pneumoniae* sequences. DFs can be readily identified by identifying unique sequences within contigs of the *Streptococcus pneumoniae* genome, such as by using well-known computer analysis software, and by generating and testing probes or amplification primers

consisting of the DF sequence in an appropriate diagnostic format which determines amplification or hybridization selectivity.

The sequences falling within the scope of the present invention are not limited to the specific sequences herein described, but also include allelic and species variations thereof. Allelic and species variations can be routinely determined by comparing the sequences provided in SEQ ID NOS:1-391, a representative fragment thereof, or a nucleotide sequence at least 95%, preferably at least 99% and most at least preferably 99.9% identical to SEQ ID NOS:1-391, with a sequence from another isolate of the same species. Furthermore, to accommodate codon variability, the invention includes nucleic acid molecules coding for the same amino acid sequences as do the specific ORFs disclosed herein. In other words, in the coding region of an ORF, substitution of one codon for another which encodes the same amino acid is expressly contemplated. Any specific sequence disclosed herein can be readily screened for errors by resequencing a particular fragment, such as an ORF, in both directions (*i.e.*, sequence both strands). Alternatively, error screening can be performed by sequencing corresponding polynucleotides of *Streptococcus pneumoniae* origin isolated by using part or all of the fragments in question as a probe or primer.

Preferred DFs of the present invention comprise at least about 17, preferably at least about 20, and more preferably at least about 50 contiguous nucleotides within an ORF set out in Tables 1-3. Most highly preferred DFs specifically hybridize to a polynucleotide containing the sequence of the ORF from which they are derived. Specific hybridization occurs even under stringent conditions defined elsewhere herein.

Each of the ORFs of the *Streptococcus pneumoniae* genome disclosed in Tables 1, 2 and 3, and the EMFs found 5' to the ORFs, can be used as polynucleotide reagents in numerous ways. For example, the sequences can be used as diagnostic probes or diagnostic amplification primers to detect the presence of a specific microbe in a sample, particularly *Streptococcus pneumoniae*. Especially preferred in this regard are ORFs such as those of Table 3, which do not match previously characterized sequences from other organisms and thus are most likely to be highly selective for *Streptococcus pneumoniae*. Also particularly preferred are ORFs that can be used to distinguish between strains of *Streptococcus pneumoniae*, particularly those that distinguish medically important strain, such as drug-resistant strains.

In addition, the fragments of the present invention, as broadly described, can be used to control gene expression through triple helix formation or antisense DNA or RNA, both of which methods are based on the binding of a polynucleotide sequence to DNA or RNA. Triple helix-formation optimally results in a shut-off of RNA transcription from DNA, while antisense RNA hybridization blocks translation of an mRNA molecule into polypeptide. Information from the sequences of the present invention can be used to design antisense and triple helix-forming oligonucleotides. Polynucleotides suitable for use in these methods are usually 20 to 40 bases in length and are designed to be complementary to a region of the gene involved in transcription, for triple-helix formation, or to the mRNA itself, for antisense inhibition. Both techniques have been demonstrated to be effective in model systems, and the requisite techniques are well known and involve routine procedures. Triple helix techniques are discussed in, for example, Lee *et al.*, *Nucl. Acids Res.* 6:3073 (1979); Cooney *et al.*, *Science* 241:456 (1988); and Dervan *et al.*, *Science* 251:1360 (1991). Antisense techniques in general are discussed in, for instance, Okano, *J. Neurochem.* 56:560 (1991) and *Oligodeoxynucleotides as Antisense Inhibitors of Gene Expression*, CRC Press, Boca Raton, FL (1988)).

The present invention further provides recombinant constructs comprising one or more fragments of the *Streptococcus pneumoniae* genomic fragments and contigs of the present invention. Certain preferred recombinant constructs of the present invention comprise a vector, such as a plasmid or viral vector, into which a fragment of the *Streptococcus pneumoniae* genome has been inserted, in a forward or reverse orientation. In the case of a vector comprising one of the ORFs of the present invention, the vector may further comprise regulatory sequences, including for example, a promoter, operably linked to the ORF. For vectors comprising the EMFs of the present invention, the vector may further comprise a marker sequence or heterologous ORF operably linked to the EMF.

Large numbers of suitable vectors and promoters are known to those of skill in the art and are commercially available for generating the recombinant constructs of the present invention. The following vectors are provided by way of example. Useful bacterial vectors include phagescript, PsiX174, pBluescript SK, pBS KS, pNH8a, pNH16a, pNH18a, pNH46a (available from Stratagene); pTrc99A, pKK223-3, pKK233-3, pDR540, pRIT5 (available from Pharmacia). Useful eukaryotic vectors include pWLneo, pSV2cat, pOG44, pXT1, pSG

(available from Stratagene) pSVK3, pBPV, pMSG, pSVL (available from Pharmacia).

Promoter regions can be selected from any desired gene using CAT (chloramphenicol transferase) vectors or other vectors with selectable markers.

- 5 Two appropriate vectors are pKK232-8 and pCM7. Particular named bacterial promoters include lacI, lacZ, T3, T7, gpt, lambda PR, and trc. Eukaryotic promoters include CMV immediate early, HSV thymidine kinase, early and late SV40, LTRs from retrovirus, and mouse metallothionein- I. Selection of the appropriate vector and promoter is well within the level of ordinary skill in the art.

- 10 The present invention further provides host cells containing any one of the isolated fragments of the *Streptococcus pneumoniae* genomic fragments and contigs of the present invention, wherein the fragment has been introduced into the host cell using known methods. The host cell can be a higher eukaryotic host cell, such as a mammalian cell, a lower eukaryotic host cell, such as a yeast cell, or  
15 a procaryotic cell, such as a bacterial cell.

- A polynucleotide of the present invention, such as a recombinant construct comprising an ORF of the present invention, may be introduced into the host by a variety of well established techniques that are standard in the art, such as calcium phosphate transfection, DEAE, dextran mediated transfection and electroporation,  
20 which are described in, for instance, Davis, L. *et al.*, BASIC METHODS IN MOLECULAR BIOLOGY (1986).

- A host cell containing one of the fragments of the *Streptococcus pneumoniae* genomic fragments and contigs of the present invention, can be used in conventional manners to produce the gene product encoded by the isolated  
25 fragment (in the case of an ORF) or can be used to produce a heterologous protein under the control of the EMF. The present invention further provides isolated polypeptides encoded by the nucleic acid fragments of the present invention or by degenerate variants of the nucleic acid fragments of the present invention. By "degenerate variant" is intended nucleotide fragments which differ  
30 from a nucleic acid fragment of the present invention (*e.g.*, an ORF) by nucleotide sequence but, due to the degeneracy of the Genetic Code, encode an identical polypeptide sequence.

Preferred nucleic acid fragments of the present invention are the ORFs and subfragments thereof depicted in Tables 2 and 3 which encode proteins.

A variety of methodologies known in the art can be utilized to obtain any one of the isolated polypeptides or proteins of the present invention. At the simplest level, the amino acid sequence can be synthesized using commercially available peptide synthesizers. This is particularly useful in producing small peptides and fragments of larger polypeptides. Such short fragments as may be obtained most readily by synthesis are useful, for example, in generating antibodies against the native polypeptide, as discussed further below.

In an alternative method, the polypeptide or protein is purified from bacterial cells which naturally produce the polypeptide or protein. One skilled in the art can readily employ well-known methods for isolating polypeptides and proteins to isolate and purify polypeptides or proteins of the present invention produced naturally by a bacterial strain, or by other methods. Methods for isolation and purification that can be employed in this regard include, but are not limited to, immunochromatography, HPLC, size-exclusion chromatography, ion-exchange chromatography, and immuno-affinity chromatography.

The polypeptides and proteins of the present invention also can be purified from cells which have been altered to express the desired polypeptide or protein. As used herein, a cell is said to be altered to express a desired polypeptide or protein when the cell, through genetic manipulation, is made to produce a polypeptide or protein which it normally does not produce or which the cell normally produces at a lower level. Those skilled in the art can readily adapt procedures for introducing and expressing either recombinant or synthetic sequences into eukaryotic or prokaryotic cells in order to generate a cell which produces one of the polypeptides or proteins of the present invention.

Any host/vector system can be used to express one or more of the ORFs of the present invention. These include, but are not limited to, eukaryotic hosts such as HeLa cells, CV-1 cell, COS cells, and Sf9 cells, as well as prokaryotic host such as *E. coli* and *B. subtilis*. The most preferred cells are those which do not normally express the particular polypeptide or protein or which expresses the polypeptide or protein at low natural level.

"Recombinant," as used herein, means that a polypeptide or protein is derived from recombinant (e.g., microbial or mammalian) expression systems. "Microbial" refers to recombinant polypeptides or proteins made in bacterial or fungal (e.g., yeast) expression systems. As a product, "recombinant microbial" defines a polypeptide or protein essentially free of native endogenous substances and unaccompanied by associated native glycosylation. Polypeptides or proteins expressed in most bacterial cultures, e.g., *E. coli*, will be free of glycosylation modifications; polypeptides or proteins expressed in yeast will have a glycosylation pattern different from that expressed in mammalian cells.

"Nucleotide sequence" refers to a heteropolymer of deoxyribonucleotides. Generally, DNA segments encoding the polypeptides and proteins provided by this invention are assembled from fragments of the *Streptococcus pneumoniae* genome and short oligonucleotide linkers, or from a series of oligonucleotides, to provide a synthetic gene which is capable of being expressed in a recombinant transcriptional unit comprising regulatory elements derived from a microbial or viral operon.

Recombinant expression vehicle or vector" refers to a plasmid or phage or virus or vector, for expressing a polypeptide from a DNA (RNA) sequence. The expression vehicle can comprise a transcriptional unit comprising an assembly of (1) a genetic regulatory elements necessary for gene expression in the host, including elements required to initiate and maintain transcription at a level sufficient for suitable expression of the desired polypeptide, including, for example, promoters and, where necessary, an enhancer and a polyadenylation signal; (2) a structural or coding sequence which is transcribed into mRNA and translated into protein, and (3) appropriate signals to initiate translation at the beginning of the desired coding region and terminate translation at its end. Structural units intended for use in yeast or eukaryotic expression systems preferably include a leader sequence enabling extracellular secretion of translated protein by a host cell. Alternatively, where recombinant protein is expressed without a leader or transport sequence, it may include an N-terminal methionine residue. This residue may or may not be subsequently cleaved from the expressed recombinant protein to provide a final product.

"Recombinant expression system" means host cells which have stably integrated a recombinant transcriptional unit into chromosomal DNA or carry the recombinant transcriptional unit extra chromosomally. The cells can be prokaryotic or eukaryotic. Recombinant expression systems as defined herein will express

heterologous polypeptides or proteins upon induction of the regulatory elements linked to the DNA segment or synthetic gene to be expressed.

Mature proteins can be expressed in mammalian cells, yeast, bacteria, or other cells under the control of appropriate promoters. Cell-free translation systems can also be employed to produce such proteins using RNAs derived from the DNA constructs of the present invention. Appropriate cloning and expression vectors for use with prokaryotic and eukaryotic hosts are described in Sambrook *et al.*, *Molecular Cloning: A Laboratory Manual*, 2<sup>nd</sup> Edition, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York (1989), the disclosure of which is hereby incorporated by reference in its entirety.

Generally, recombinant expression vectors will include origins of replication and selectable markers permitting transformation of the host cell, *e.g.*, the ampicillin resistance gene of *E. coli* and *S. cerevisiae* TRP1 gene, and a promoter derived from a highly expressed gene to direct transcription of a downstream structural sequence. Such promoters can be derived from operons encoding glycolytic enzymes such as 3-phosphoglycerate kinase (PGK), alpha-factor, acid phosphatase, or heat shock proteins, among others. The heterologous structural sequence is assembled in appropriate phase with translation initiation and termination sequences, and preferably, a leader sequence capable of directing secretion of translated protein into the periplasmic space or extracellular medium. Optionally, the heterologous sequence can encode a fusion protein including an N-terminal identification peptide imparting desired characteristics, *e.g.*, stabilization or simplified purification of expressed recombinant product.

Useful expression vectors for bacterial use are constructed by inserting a structural DNA sequence encoding a desired protein together with suitable translation initiation and termination signals in operable reading phase with a functional promoter. The vector will comprise one or more phenotypic selectable markers and an origin of replication to ensure maintenance of the vector and, when desirable, provide amplification within the host.

Suitable prokaryotic hosts for transformation include strains of *E. coli*, *B. subtilis*, *Salmonella typhimurium* and various species within the genera *Pseudomonas* and *Streptomyces*. Others may, also be employed as a matter of choice.

As a representative but non-limiting example, useful expression vectors for bacterial use can comprise a selectable marker and bacterial origin of replication

derived from commercially available plasmids comprising genetic elements of the well known cloning vector pBR322 (ATCC 37017). Such commercial vectors include, for example, pKK223-3 (available from Pharmacia Fine Chemicals, Uppsala, Sweden) and GEM 1 (available from Promega Biotec, Madison, WI, USA). These pBR322 "backbone" sections are combined with an appropriate promoter and the structural sequence to be expressed.

Following transformation of a suitable host strain and growth of the host strain to an appropriate cell density, the selected promoter, where it is inducible, is derepressed or induced by appropriate means (*e.g.*, temperature shift or chemical induction) and cells are cultured for an additional period to provide for expression of the induced gene product. Thereafter cells are typically harvested, generally by centrifugation, disrupted to release expressed protein, generally by physical or chemical means, and the resulting crude extract is retained for further purification.

Various mammalian cell culture systems can also be employed to express recombinant protein. Examples of mammalian expression systems include the COS-7 lines of monkey kidney fibroblasts, described in Gluzman, *Cell* 23:175 (1981), and other cell lines capable of expressing a compatible vector, for example, the C127, 3T3, CHO, HeLa and BHK cell lines.

Mammalian expression vectors will comprise an origin of replication, a suitable promoter and enhancer, and also any necessary ribosome binding sites, polyadenylation site, splice donor and acceptor sites, transcriptional termination sequences, and 5' flanking nontranscribed sequences. DNA sequences derived from the SV40 viral genome, for example, SV40 origin, early promoter, enhancer, splice, and polyadenylation sites may be used to provide the required nontranscribed genetic elements.

Recombinant polypeptides and proteins produced in bacterial culture is usually isolated by initial extraction from cell pellets, followed by one or more salting-out, aqueous ion exchange or size exclusion chromatography steps. Microbial cells employed in expression of proteins can be disrupted by any convenient method, including freeze-thaw cycling, sonication, mechanical disruption, or use of cell lysing agents. Protein refolding steps can be used, as necessary, in completing configuration of the mature protein. Finally, high performance liquid chromatography (HPLC) can be employed for final purification steps.



The present invention further includes isolated polypeptides, proteins and nucleic acid molecules which are substantially equivalent to those herein described. As used herein, substantially equivalent can refer both to nucleic acid and amino acid sequences, for example a mutant sequence, that varies from a reference sequence by one or more substitutions, deletions, or additions, the net effect of which does not result in an adverse functional dissimilarity between reference and subject sequences. For purposes of the present invention, sequences having equivalent biological activity, and equivalent expression characteristics are considered substantially equivalent. For purposes of determining equivalence, truncation of the mature sequence should be disregarded.

The invention further provides methods of obtaining homologs from other strains of *Streptococcus pneumoniae*, of the fragments of the *Streptococcus pneumoniae* genome of the present invention and homologs of the proteins encoded by the ORFs of the present invention. As used herein, a sequence or protein of *Streptococcus pneumoniae* is defined as a homolog of a fragment of the *Streptococcus pneumoniae* fragments or contigs or a protein encoded by one of the ORFs of the present invention, if it shares significant homology to one of the fragments of the *Streptococcus pneumoniae* genome of the present invention or a protein encoded by one of the ORFs of the present invention. Specifically, by using the sequence disclosed herein as a probe or as primers, and techniques such as PCR cloning and colony/plaque hybridization, one skilled in the art can obtain homologs.

As used herein, two nucleic acid molecules or proteins are said to "share significant homology" if the two contain regions which possess greater than 85% sequence (amino acid or nucleic acid) homology. Preferred homologs in this regard are those with more than 90% homology. Especially preferred are those with 93% or more homology. Among especially preferred homologs those with 95% or more homology are particularly preferred. Very particularly preferred among these are those with 97% and even more particularly preferred among these are homologs with 99% or more homology. The most preferred homologs among these are those with 99.9% homology or more. It will be understood that, among measures of homology, identity is particularly preferred in this regard.

Region specific primers or probes derived from the nucleotide sequence provided in SEQ ID NOS:1-391 or from a nucleotide sequence at least 95%, particularly at least 99%, especially at least 99.5% identical to a sequence of SEQ

ID NOS:1-391 can be used to prime DNA synthesis and PCR amplification, as well as to identify colonies containing cloned DNA encoding a homolog. Methods suitable to this aspect of the present invention are well known and have been described in great detail in many publications such as, for example, Innis *et al.*,  
5 *PCR Protocols*, Academic Press, San Diego, CA (1990).

When using primers derived from SEQ ID NOS:1-391 or from a nucleotide sequence having an aforementioned identity to a sequence of SEQ ID NOS:1-391, one skilled in the art will recognize that by employing high stringency conditions (*e.g.*, annealing at 50-60°C in 6X SSPC and 50% formamide, and washing at 50-  
10 65°C in 0.5X SSPC) only sequences which are greater than 75% homologous to the primer will be amplified. By employing lower stringency conditions (*e.g.*, hybridizing at 35-37°C in 5X SSPC and 40-45% formamide, and washing at 42°C in 0.5X SSPC), sequences which are greater than 40-50% homologous to the primer will also be amplified.

When using DNA probes derived from SEQ ID NOS:1-391, or from a nucleotide sequence having an aforementioned identity to a sequence of SEQ ID NOS:1-391, for colony/plaque hybridization, one skilled in the art will recognize that by employing high stringency conditions (*e.g.*, hybridizing at 50- 65°C in 5X  
15 SSPC and 50% formamide, and washing at 50- 65°C in 0.5X SSPC), sequences having regions which are greater than 90% homologous to the probe can be obtained, and that by employing lower stringency conditions (*e.g.*, hybridizing at 35-37°C in 5X SSPC and 40-45% formamide, and washing at 42°C in 0.5X SSPC), sequences having regions which are greater than 35-45% homologous to the probe will be obtained.

Any organism can be used as the source for homologs of the present invention so long as the organism naturally expresses such a protein or contains genes encoding the same. The most preferred organism for isolating homologs are bacteria which are closely related to *Streptococcus pneumoniae*.

## 30 ILLUSTRATIVE USES OF COMPOSITIONS OF THE INVENTION

Each ORF provided in Tables 1 and 2 is identified with a function by homology to a known gene or polypeptide. As a result, one skilled in the art can use the polypeptides of the present invention for commercial, therapeutic and  
35 industrial purposes consistent with the type of putative identification of the

polypeptide. Such identifications permit one skilled in the art to use the *Streptococcus pneumoniae* ORFs in a manner similar to the known type of sequences for which the identification is made; for example, to ferment a particular sugar source or to produce a particular metabolite. A variety of reviews illustrative of this aspect of the invention are available, including the following reviews on the industrial use of enzymes, for example, BIOCHEMICAL ENGINEERING AND BIOTECHNOLOGY HANDBOOK, 2nd Ed., MacMillan Publications, Ltd. NY (1991) and BIOCATALYSTS IN ORGANIC SYNTHESSES, Tramper *et al.*, Eds., Elsevier Science Publishers, Amsterdam, The Netherlands (1985). A variety of exemplary uses that illustrate this and similar aspects of the present invention are discussed below.

### 1. Biosynthetic Enzymes

Open reading frames encoding proteins involved in mediating the catalytic reactions involved in intermediary and macromolecular metabolism, the biosynthesis of small molecules, cellular processes and other functions includes enzymes involved in the degradation of the intermediary products of metabolism, enzymes involved in central intermediary metabolism, enzymes involved in respiration, both aerobic and anaerobic, enzymes involved in fermentation, enzymes involved in ATP proton motor force conversion, enzymes involved in broad regulatory function, enzymes involved in amino acid synthesis, enzymes involved in nucleotide synthesis, enzymes involved in cofactor and vitamin synthesis, can be used for industrial biosynthesis.

The various metabolic pathways present in *Streptococcus pneumoniae* can be identified based on absolute nutritional requirements as well as by examining the various enzymes identified in Table 1-3 and SEQ ID NOS:1-391.

Of particular interest are polypeptides involved in the degradation of intermediary metabolites as well as non-macromolecular metabolism. Such enzymes include amylases, glucose oxidases, and catalase.

Proteolytic enzymes are another class of commercially important enzymes. Proteolytic enzymes find use in a number of industrial processes including the processing of flax and other vegetable fibers, in the extraction, clarification and depectinization of fruit juices, in the extraction of vegetables' oil and in the maceration of fruits and vegetables to give unicellular fruits. A detailed review of the proteolytic enzymes used in the food industry is provided in Rombouts *et al.*,

*Symbiosis* 21:79 (1986) and Voragen *et al.* in *Biocatalysts In Agricultural Biotechnology*, Whitaker *et al.*, Eds., *American Chemical Society Symposium Series* 389:93 (1989).

The metabolism of sugars is an important aspect of the primary metabolism of *Streptococcus pneumoniae*. Enzymes involved in the degradation of sugars, such as, particularly, glucose, galactose, fructose and xylose, can be used in industrial fermentation. Some of the important sugar transforming enzymes, from a commercial viewpoint, include sugar isomerases such as glucose isomerase. Other metabolic enzymes have found commercial use such as glucose oxidases which produces ketogulonic acid (KGA). KGA is an intermediate in the commercial production of ascorbic acid using the Reichstein's procedure, as described in Krueger *et al.*, *Biotechnology* 6(A). Rhine *et al.*, Eds., Verlag Press, Weinheim, Germany (1984).

Glucose oxidase (GOD) is commercially available and has been used in purified form as well as in an immobilized form for the deoxygenation of beer. See, for instance, Hartmeir *et al.*, *Biotechnology Letters* 1:21 (1979). The most important application of GOD is the industrial scale fermentation of gluconic acid. Market for gluconic acids which are used in the detergent, textile, leather, photographic, pharmaceutical, food, feed and concrete industry, as described, for example, in Bigelis *et al.*, beginning on page 357 in *GENE MANIPULATIONS AND FUNGI*; Benett *et al.*, Eds., Academic Press, New York (1985). In addition to industrial applications, GOD has found applications in medicine for quantitative determination of glucose in body fluids recently in biotechnology for analyzing syrups from starch and cellulose hydrosylates. This application is described in Owusu *et al.*, *Biochem. et Biophysica. Acta.* 872:83 (1986), for instance.

The main sweetener used in the world today is sugar which comes from sugar beets and sugar cane. In the field of industrial enzymes, the glucose isomerase process shows the largest expansion in the market today. Initially, soluble enzymes were used and later immobilized enzymes were developed (Krueger *et al.*, *Biotechnology, The Textbook of Industrial Microbiology*, Sinauer Associated Incorporated, Sunderland, Massachusetts (1990)). Today, the use of glucose- produced high fructose syrups is by far the largest industrial business using immobilized enzymes. A review of the industrial use of these enzymes is provided by Jorgensen, *Starch* 40:307 (1988).

Proteinases, such as alkaline serine proteinases, are used as detergent additives and thus represent one of the largest volumes of microbial enzymes used in the industrial sector. Because of their industrial importance, there is a large body of published and unpublished information regarding the use of these enzymes in industrial processes. (See Faultman *et al.*, Acid Proteases Structure Function and Biology, Tang, J., ed., Plenum Press, New York (1977) and Godfrey *et al.*, Industrial Enzymes, MacMillan Publishers, Surrey, UK (1983) and Hepner *et al.*, Report Industrial Enzymes by 1990, Hel Hepner & Associates, London (1986)).

Another class of commercially usable proteins of the present invention are the microbial lipases, described by, for instance, Macrae *et al.*, *Philosophical Transactions of the Chiral Society of London* 310:227 (1985) and Poserke, *Journal of the American Oil Chemist Society* 61:1758 (1984). A major use of lipases is in the fat and oil industry for the production of neutral glycerides using lipase catalyzed inter-esterification of readily available triglycerides. Application of lipases include the use as a detergent additive to facilitate the removal of fats from fabrics in the course of the washing procedures.

The use of enzymes, and in particular microbial enzymes, as catalyst for key steps in the synthesis of complex organic molecules is gaining popularity at a great rate. One area of great interest is the preparation of chiral intermediates. Preparation of chiral intermediates is of interest to a wide range of synthetic chemists particularly those scientists involved with the preparation of new pharmaceuticals, agrochemicals, fragrances and flavors. (See Davies *et al.*, *Recent Advances in the Generation of Chiral Intermediates Using Enzymes*, CRC Press, Boca Raton, Florida (1990)). The following reactions catalyzed by enzymes are of interest to organic chemists: hydrolysis of carboxylic acid esters, phosphate esters, amides and nitriles, esterification reactions, trans-esterification reactions, synthesis of amides, reduction of alkanones and oxoalkanates, oxidation of alcohols to carbonyl compounds, oxidation of sulfides to sulfoxides, and carbon bond forming reactions such as the aldol reaction.

When considering the use of an enzyme encoded by one of the ORFs of the present invention for biotransformation and organic synthesis it is sometimes necessary to consider the respective advantages and disadvantages of using a microorganism as opposed to an isolated enzyme. Pros and cons of using a whole cell system on the one hand or an isolated partially purified enzyme on the other

hand, has been described in detail by Bud *et al.*, Chemistry in Britain (1987), p. 127.

Amino transferases, enzymes involved in the biosynthesis and metabolism of amino acids, are useful in the catalytic production of amino acids. The advantages of using microbial based enzyme systems is that the amino transferase enzymes catalyze the stereo- selective synthesis of only L-amino acids and generally possess uniformly high catalytic rates. A description of the use of amino transferases for amino acid production is provided by Roselle-David, *Methods of Enzymology* 136:479 (1987).

Another category of useful proteins encoded by the ORFs of the present invention include enzymes involved in nucleic acid synthesis, repair, and recombination.

## 2. Generation of Antibodies

As described here, the proteins of the present invention, as well as homologs thereof, can be used in a variety of procedures and methods known in the art which are currently applied to other proteins. The proteins of the present invention can further be used to generate an antibody which selectively binds the protein. Such antibodies can be either monoclonal or polyclonal antibodies, as well fragments of these antibodies, and humanized forms.

The invention further provides antibodies which selectively bind to one of the proteins of the present invention and hybridomas which produce these antibodies. A hybridoma is an immortalized cell line which is capable of secreting a specific monoclonal antibody.

In general, techniques for preparing polyclonal and monoclonal antibodies as well as hybridomas capable of producing the desired antibody are well known in the art (Campbell, A. M., *Monoclonal Antibody Technology: Laboratory Techniques In Biochemistry And Molecular Biology*, Elsevier Science Publishers, Amsterdam, The Netherlands (1984); St. Groth *et al.*, *J. Immunol. Methods* 35: 1-21 (1980), Kohler and Milstein, *Nature* 256:495-497 (1975)), the trioma technique, the human B-cell hybridoma technique (Kozbor *et al.*, *Immunology Today* 4:72 (1983), pgs. 77-96 of Cole *et al.*, in *Monoclonal Antibodies And Cancer Therapy*, Alan R. Liss, Inc. (1985)). Any animal (mouse, rabbit, etc.) which is known to produce antibodies can be immunized with the pseudogene polypeptide. Methods for immunization are well known in the art. Such methods

include subcutaneous or interperitoneal injection of the polypeptide. One skilled in the art will recognize that the amount of the protein encoded by the ORF of the present invention used for immunization will vary based on the animal which is immunized, the antigenicity of the peptide and the site of injection.

5       The protein which is used as an immunogen may be modified or administered in an adjuvant in order to increase the protein's antigenicity. Methods of increasing the antigenicity of a protein are well known in the art and include, but are not limited to coupling the antigen with a heterologous protein (such as globulin or galactosidase) or through the inclusion of an adjuvant during immunization.

10       For monoclonal antibodies, spleen cells from the immunized animals are removed, fused with myeloma cells, such as SP2/0-Ag14 myeloma cells, and allowed to become monoclonal antibody producing hybridoma cells.

Any one of a number of methods well known in the art can be used to identify the hybridoma cell which produces an antibody with the desired characteristics. These include screening the hybridomas with an ELISA assay, western blot analysis, or radioimmunoassay (Lutz *et al.*, *Exp. Cell Res.* 175:109-124 (1988)).

Hybridomas secreting the desired antibodies are cloned and the class and subclass is determined using procedures known in the art (Campbell, A. M., *Monoclonal Antibody Technology: Laboratory Techniques in Biochemistry and Molecular Biology*, Elsevier Science Publishers, Amsterdam, The Netherlands (1984)).

Techniques described for the production of single chain antibodies (U. S. Patent 4,946,778) can be adapted to produce single chain antibodies to proteins of the present invention.

For polyclonal antibodies, antibody containing antisera is isolated from the immunized animal and is screened for the presence of antibodies with the desired specificity using one of the above-described procedures.

The present invention further provides the above- described antibodies in detectably labelled form. Antibodies can be detectably labelled through the use of radioisotopes, affinity labels (such as biotin, avidin, *etc.*), enzymatic labels (such as horseradish peroxidase, alkaline phosphatase, *etc.*) fluorescent labels (such as FITC or rhodamine, *etc.*), paramagnetic atoms, *etc.* Procedures for accomplishing such labeling are well-known in the art, for example see Sternberger *et al.*, *J. Histochem. Cytochem.* 18:315 (1970); Bayer, E. A. *et al.*, *Meth. Enzym.* 62:308

(1979); Engval, E. *et al.*, *Immunol.* 109:129 (1972); Goding, J. W., *J. Immunol. Meth.* 13:215 (1976)).

The labeled antibodies of the present invention can be used for *in vitro*, *in vivo*, and *in situ* assays to identify cells or tissues in which a fragment of the  
5 *Streptococcus pneumoniae* genome is expressed.

The present invention further provides the above-described antibodies immobilized on a solid support. Examples of such solid supports include plastics such as polycarbonate, complex carbohydrates such as agarose and sepharose, acrylic resins and such as polyacrylamide and latex beads. Techniques for  
10 coupling antibodies to such solid supports are well known in the art (Weir, D. M. *et al.*, "Handbook of Experimental Immunology" 4th Ed., Blackwell Scientific Publications, Oxford, England, Chapter 10 (1986); Jacoby, W. D. *et al.*, *Meth. Enzym.* 34 Academic Press, N. Y. (1974)). The immobilized antibodies of the present invention can be used for *in vitro*, *in vivo*, and *in situ* assays as well as for  
15 immunoaffinity purification of the proteins of the present invention.

### 3. Diagnostic Assays and Kits

The present invention further provides methods to identify the expression of one of the ORFs of the present invention, or homolog thereof, in a test sample,  
20 using one of the DFs or antibodies of the present invention.

In detail, such methods comprise incubating a test sample with one or more of the antibodies or one or more of the DFs of the present invention and assaying for binding of the DFs or antibodies to components within the test sample.

Conditions for incubating a DF or antibody with a test sample vary.  
25 Incubation conditions depend on the format employed in the assay, the detection methods employed, and the type and nature of the DF or antibody used in the assay. One skilled in the art will recognize that any one of the commonly available hybridization, amplification or immunological assay formats can readily be adapted to employ the DFs or antibodies of the present invention. Examples of such assays  
30 can be found in Chard, T., *An Introduction to Radioimmunoassay and Related Techniques*, Elsevier Science Publishers, Amsterdam, The Netherlands (1986); Bullock, G. R. *et al.*, *Techniques in Immunocytochemistry*, Academic Press, Orlando, FL Vol. 1 (1982), Vol. 2 (1983), Vol. 3 (1985); Tijssen, P., *Practice and Theory of Enzyme Immunoassays: Laboratory Techniques in Biochemistry and*



*Molecular Biology*, Elsevier Science Publishers, Amsterdam, The Netherlands (1985).

The test samples of the present invention include cells, protein or membrane extracts of cells, or biological fluids such as sputum, blood, serum, plasma, or urine. The test sample used in the above-described method will vary based on the assay format, nature of the detection method and the tissues, cells or extracts used as the sample to be assayed. Methods for preparing protein extracts or membrane extracts of cells are well known in the art and can be readily be adapted in order to obtain a sample which is compatible with the system utilized.

In another embodiment of the present invention, kits are provided which contain the necessary reagents to carry out the assays of the present invention.

Specifically, the invention provides a compartmentalized kit to receive, in close confinement, one or more containers which comprises: (a) a first container comprising one of the DFs or antibodies of the present invention; and (b) one or more other containers comprising one or more of the following: wash reagents, reagents capable of detecting presence of a bound DF or antibody.

In detail, a compartmentalized kit includes any kit in which reagents are contained in separate containers. Such containers include small glass containers, plastic containers or strips of plastic or paper. Such containers allows one to efficiently transfer reagents from one compartment to another compartment such that the samples and reagents are not cross-contaminated, and the agents or solutions of each container can be added in a quantitative fashion from one compartment to another. Such containers will include a container which will accept the test sample, a container which contains the antibodies used in the assay, containers which contain wash reagents (such as phosphate buffered saline, Tris-buffers, *etc.*), and containers which contain the reagents used to detect the bound antibody or DF.

Types of detection reagents include labelled nucleic acid probes, labelled secondary antibodies, or in the alternative, if the primary antibody is labelled, the enzymatic, or antibody binding reagents which are capable of reacting with the labelled antibody. One skilled in the art will readily recognize that the disclosed DFs and antibodies of the present invention can be readily incorporated into one of the established kit formats which are well known in the art.

#### 4. Screening Assay for Binding Agents

Using the isolated proteins of the present invention, the present invention further provides methods of obtaining and identifying agents which bind to a protein encoded by one of the ORFs of the present invention or to one of the fragments and the *Streptococcus pneumoniae* fragment and contigs herein  
5 described.

In general, such methods comprise steps of:

- (a) contacting an agent with an isolated protein encoded by one of the ORFs of the present invention, or an isolated fragment of the *Streptococcus pneumoniae* genome; and
- 10 (b) determining whether the agent binds to said protein or said fragment.

The agents screened in the above assay can be, but are not limited to, peptides, carbohydrates, vitamin derivatives, or other pharmaceutical agents. The agents can be selected and screened at random or rationally selected or designed using protein modeling techniques.

- 15 For random screening, agents such as peptides, carbohydrates, pharmaceutical agents and the like are selected at random and are assayed for their ability to bind to the protein encoded by the ORF of the present invention.

Alternatively, agents may be rationally selected or designed. As used herein, an agent is said to be "rationally selected or designed" when the agent is  
20 chosen based on the configuration of the particular protein. For example, one skilled in the art can readily adapt currently available procedures to generate peptides, pharmaceutical agents and the like capable of binding to a specific peptide sequence in order to generate rationally designed antipeptide peptides, for example see Hurby *et al.*, "Application of Synthetic Peptides: Antisense Peptides," in  
25 *Synthetic Peptides, A User's Guide*, W. H. Freeman, NY (1992), pp. 289-307, and Kasieczak *et al.*, *Biochemistry* 28:9230-8 (1989), or pharmaceutical agents, or the like.

In addition to the foregoing, one class of agents of the present invention, as broadly described, can be used to control gene expression through binding to one  
30 of the ORFs or EMFs of the present invention. As described above, such agents can be randomly screened or rationally designed/selected. Targeting the ORF or EMF allows a skilled artisan to design sequence specific or element specific agents, modulating the expression of either a single ORF or multiple ORFs which rely on the same EMF for expression control.

One class of DNA binding agents are agents which contain base residues which hybridize or form a triple helix by binding to DNA or RNA. Such agents can be based on the classic phosphodiester, ribonucleic acid backbone, or can be a variety of sulfhydryl or polymeric derivatives which have base attachment capacity.

- 5 Agents suitable for use in these methods usually contain 20 to 40 bases and are designed to be complementary to a region of the gene involved in transcription (triple helix - see Lee *et al.*, *Nucl. Acids Res.* 6:3073 (1979); Cooney *et al.*, *Science* 241:456 (1988); and Dervan *et al.*, *Science* 251:1360 (1991)) or to the mRNA itself (antisense - Okano, *J. Neurochem.* 56:560 (1991);
- 10 *Oligodeoxynucleotides as Antisense Inhibitors of Gene Expression*, CRC Press, Boca Raton, FL (1988)). Triple helix- formation optimally results in a shut-off of RNA transcription from DNA, while antisense RNA hybridization blocks translation of an mRNA molecule into polypeptide. Both techniques have been demonstrated to be effective in model systems. Information contained in the
- 15 sequences of the present invention can be used to design antisense and triple helix-forming oligonucleotides, and other DNA binding agents.

## 5. Pharmaceutical Compositions and Vaccines

- The present invention further provides pharmaceutical agents which can be
- 20 used to modulate the growth or pathogenicity of *Streptococcus pneumoniae*, or another related organism, *in vivo* or *in vitro*. As used herein, a "pharmaceutical agent" is defined as a composition of matter which can be formulated using known techniques to provide a pharmaceutical compositions. As used herein, the "pharmaceutical agents of the present invention" refers the pharmaceutical agents
- 25 which are derived from the proteins encoded by the ORFs of the present invention or are agents which are identified using the herein described assays.

- As used herein, a pharmaceutical agent is said to "modulate the growth pathogenicity of *Streptococcus pneumoniae* or a related organism, *in vivo* or *in vitro*," when the agent reduces the rate of growth, rate of division, or viability of
- 30 the organism in question. The pharmaceutical agents of the present invention can modulate the growth or pathogenicity of an organism in many fashions, although an understanding of the underlying mechanism of action is not needed to practice the use of the pharmaceutical agents of the present invention. Some agents will modulate the growth by binding to an important protein thus blocking the biological
- 35 activity of the protein, while other agents may bind to a component of the outer

surface of the organism blocking attachment or rendering the organism more prone to act the bodies nature immune system. Alternatively, the agent may comprise a protein encoded by one of the ORFs of the present invention and serve as a vaccine. The development and use of a vaccine based on outer membrane  
5 components are well known in the art.

As used herein, a "related organism" is a broad term which refers to any organism whose growth can be modulated by one of the pharmaceutical agents of the present invention. In general, such an organism will contain a homolog of the protein which is the target of the pharmaceutical agent or the protein used as a  
10 vaccine. As such, related organisms do not need to be bacterial but may be fungal or viral pathogens.

The pharmaceutical agents and compositions of the present invention may be administered in a convenient manner, such as by the oral, topical, intravenous, intraperitoneal, intramuscular, subcutaneous, intranasal or intradermal routes. The  
15 pharmaceutical compositions are administered in an amount which is effective for treating and/or prophylaxis of the specific indication. In general, they are administered in an amount of at least about 1 mg/kg body weight and in most cases they will be administered in an amount not in excess of about 1 g/kg body weight per day. In most cases, the dosage is from about 0.1 mg/kg to about 10 g/kg body  
20 weight daily, taking into account the routes of administration, symptoms, *etc.*

The agents of the present invention can be used in native form or can be modified to form a chemical derivative. As used herein, a molecule is said to be a "chemical derivative" of another molecule when it contains additional chemical  
25 moieties not normally a part of the molecule. Such moieties may improve the molecule's solubility, absorption, biological half life, *etc.* The moieties may alternatively decrease the toxicity of the molecule, eliminate or attenuate any undesirable side effect of the molecule, *etc.* Moieties capable of mediating such effects are disclosed in, among other sources, REMINGTON'S PHARMACEUTICAL SCIENCES (1980) cited elsewhere herein.

For example, such moieties may change an immunological character of the functional derivative, such as affinity for a given antibody. Such changes in immunomodulation activity are measured by the appropriate assay, such as a  
35 competitive type immunoassay. Modifications of such protein properties as redox or thermal stability, biological half-life, hydrophobicity, susceptibility to proteolytic degradation or the tendency to aggregate with carriers or into multimers also may

be effected in this way and can be assayed by methods well known to the skilled artisan.

The therapeutic effects of the agents of the present invention may be obtained by providing the agent to a patient by any suitable means (*e.g.*, inhalation, intravenously, intramuscularly, subcutaneously, enterally, or parenterally). It is preferred to administer the agent of the present invention so as to achieve an effective concentration within the blood or tissue in which the growth of the organism is to be controlled. To achieve an effective blood concentration, the preferred method is to administer the agent by injection. The administration may be by continuous infusion, or by single or multiple injections.

In providing a patient with one of the agents of the present invention, the dosage of the administered agent will vary depending upon such factors as the patient's age, weight, height, sex, general medical condition, previous medical history, *etc.* In general, it is desirable to provide the recipient with a dosage of agent which is in the range of from about 1 pg/kg to 10 mg/kg (body weight of patient), although a lower or higher dosage may be administered. The therapeutically effective dose can be lowered by using combinations of the agents of the present invention or another agent.

As used herein, two or more compounds or agents are said to be administered "in combination" with each other when either (1) the physiological effects of each compound, or (2) the serum concentrations of each compound can be measured at the same time. The composition of the present invention can be administered concurrently with, prior to, or following the administration of the other agent.

The agents of the present invention are intended to be provided to recipient subjects in an amount sufficient to decrease the rate of growth (as defined above) of the target organism.

The administration of the agent(s) of the invention may be for either a "prophylactic" or "therapeutic" purpose. When provided prophylactically, the agent(s) are provided in advance of any symptoms indicative of the organisms growth. The prophylactic administration of the agent(s) serves to prevent, attenuate, or decrease the rate of onset of any subsequent infection. When provided therapeutically, the agent(s) are provided at (or shortly after) the onset of an indication of infection. The therapeutic administration of the compound(s)

serves to attenuate the pathological symptoms of the infection and to increase the rate of recovery.

The agents of the present invention are administered to a subject, such as a mammal, or a patient, in a pharmaceutically acceptable form and in a therapeutically effective concentration. A composition is said to be "pharmacologically acceptable" if its administration can be tolerated by a recipient patient. Such an agent is said to be administered in a "therapeutically effective amount" if the amount administered is physiologically significant. An agent is physiologically significant if its presence results in a detectable change in the physiology of a recipient patient.

The agents of the present invention can be formulated according to known methods to prepare pharmaceutically useful compositions, whereby these materials, or their functional derivatives, are combined in a mixture with a pharmaceutically acceptable carrier vehicle. Suitable vehicles and their formulation, inclusive of other human proteins, *e.g.*, human serum albumin, are described, for example, in REMINGTON'S PHARMACEUTICAL SCIENCES, 16<sup>th</sup> Ed., Osol, A., Ed., Mack Publishing, Easton PA (1980). In order to form a pharmaceutically acceptable composition suitable for effective administration, such compositions will contain an effective amount of one or more of the agents of the present invention, together with a suitable amount of carrier vehicle.

Additional pharmaceutical methods may be employed to control the duration of action. Control release preparations may be achieved through the use of polymers to complex or absorb one or more of the agents of the present invention. The controlled delivery may be effectuated by a variety of well known techniques, including formulation with macromolecules such as, for example, polyesters, polyamino acids, polyvinyl, pyrrolidone, ethylenevinylacetate, methylcellulose, carboxymethylcellulose, or protamine, sulfate, adjusting the concentration of the macromolecules and the agent in the formulation, and by appropriate use of methods of incorporation, which can be manipulated to effectuate a desired time course of release. Another possible method to control the duration of action by controlled release preparations is to incorporate agents of the present invention into particles of a polymeric material such as polyesters, polyamino acids, hydrogels, poly(lactic acid) or ethylene vinylacetate copolymers. Alternatively, instead of incorporating these agents into polymeric particles, it is possible to entrap these materials in microcapsules prepared, for example, by coacervation techniques or by interfacial polymerization with, for example, hydroxymethylcellulose or gelatine-

microcapsules and poly(methylmethacrylate) microcapsules, respectively, or in colloidal drug delivery systems, for example, liposomes, albumin microspheres, microemulsions, nanoparticles, and nanocapsules or in macroemulsions. Such techniques are disclosed in REMINGTON'S PHARMACEUTICAL SCIENCES  
5 (1980).

The invention further provides a pharmaceutical pack or kit comprising one or more containers filled with one or more of the ingredients of the pharmaceutical compositions of the invention. Associated with such container(s) can be a notice in the form prescribed by a governmental agency regulating the manufacture, use or  
10 sale of pharmaceuticals or biological products, which notice reflects approval by the agency of manufacture, use or sale for human administration.

In addition, the agents of the present invention may be employed in conjunction with other therapeutic compounds.

## 15 6. Shot-Gun Approach to Megabase DNA Sequencing

The present invention further demonstrates that a large sequence can be sequenced using a random shotgun approach. This procedure, described in detail in the examples that follow, has eliminated the up front cost of isolating and ordering overlapping or contiguous subclones prior to the start of the sequencing  
20 protocols.

Certain aspects of the present invention are described in greater detail in the examples that follow. The examples are provided by way of illustration. Other aspects and embodiments of the present invention are contemplated by the inventors, as will be clear to those of skill in the art from reading the present  
25 disclosure.

## **ILLUSTRATIVE EXAMPLES**

### **LIBRARIES AND SEQUENCING**

#### 30 1. Shotgun Sequencing Probability Analysis

The overall strategy for a shotgun approach to whole genome sequencing follows from the Lander and Waterman (Landerman and Waterman, *Genomics* 2:231 (1988)) application of the equation for the Poisson distribution. According to this treatment, the probability,  $P$ , that any given base in a sequence of size  $L$ , in  
35 nucleotides, is not sequenced after a certain amount,  $n$ , in nucleotides, of random  
0

sequence has been determined can be calculated by the equation  $P = e^{-m}$ , where  $m$  is  $L/n$ , the fold coverage. For instance, for a genome of 2.8 Mb,  $m=1$  when 2.8 Mb of sequence has been randomly generated (1X coverage). At that point,  $P = e^{-1} = 0.37$ . The probability that any given base has not been sequenced is the same as the probability that any region of the whole sequence  $L$  has not been determined and, therefore, is equivalent to the fraction of the whole sequence that has yet to be determined. Thus, at one-fold coverage, approximately 37% of a polynucleotide of size  $L$ , in nucleotides has not been sequenced. When 14 Mb of sequence has been generated, coverage is 5X for a 2.8 Mb and the unsequenced fraction drops to .0067 or 0.67%. 5X coverage of a 2.8 Mb sequence can be attained by sequencing approximately 17,000 random clones from both insert ends with an average sequence read length of 410 bp.

Similarly, the total gap length,  $G$ , is determined by the equation  $G = Le^{-m}$ , and the average gap size,  $g$ , follows the equation,  $g = L/n$ . Thus, 5X coverage leaves about 240 gaps averaging about 82 bp in size in a sequence of a polynucleotide 2.8 Mb long.

The treatment above is essentially that of Lander and Waterman, *Genomics* 2: 231 (1988).

## 2. Random Library Construction

In order to approximate the random model described above during actual sequencing, a nearly ideal library of cloned genomic fragments is required. The following library construction procedure was developed to achieve this end.

*Streptococcus pneumoniae* DNA is prepared by phenol extraction. A mixture containing 200 µg DNA in 1.0 ml of 300 mM sodium acetate, 10 mM Tris-HCl, 1 mM Na-EDTA, 50% glycerol is processed through a nebulizer (PI Medical Products) with a stream of nitrogen adjusted to 35 Kpa for 2 minutes. The sonicated DNA is ethanol precipitated and redissolved in 500 µl TE buffer.

To create blunt-ends, a 100 µl aliquot of the resuspended DNA is digested with 5 units of BAL31 nuclease (New England BioLabs) for 10 min at 30°C in 200 µl BAL31 buffer. The digested DNA is phenol-extracted, ethanol-precipitated, redissolved in 100 µl TE buffer, and then size-fractionated by electrophoresis through a 1.0% low melting temperature agarose gel. The section containing DNA fragments 1.6-2.0 kb in size is excised from the gel, and the LGT agarose is melted and the resulting solution is extracted with phenol to separate the agarose from the



DNA. DNA is ethanol precipitated and redissolved in 20  $\mu$ l of TE buffer for ligation to vector.

A two-step ligation procedure is used to produce a plasmid library with 97% inserts, of which >99% were single inserts. The first ligation mixture (50  $\mu$ l) contains 2  $\mu$ g of DNA fragments, 2  $\mu$ g pUC18 DNA (Pharmacia) cut with SmaI and dephosphorylated with bacterial alkaline phosphatase, and 10 units of T4 ligase (GIBCO/BRL) and is incubated at 14°C for 4 hr. The ligation mixture then is phenol extracted and ethanol precipitated, and the precipitated DNA is dissolved in 20  $\mu$ l TE buffer and electrophoresed on a 1.0% low melting agarose gel. Discrete bands in a ladder are visualized by ethidium bromide-staining and UV illumination and identified by size as insert (I), vector (v), v+I, v+2i, v+3i, etc. The portion of the gel containing v+I DNA is excised and the v+I DNA is recovered and resuspended into 20  $\mu$ l TE. The v+I DNA then is blunt-ended by T4 polymerase treatment for 5 min. at 37°C in a reaction mixture (50  $\mu$ l) containing the v+I linears, 500  $\mu$ M each of the 4 dNTPs, and 9 units of T4 polymerase (New England BioLabs), under recommended buffer conditions. After phenol extraction and ethanol precipitation the repaired v+I linears are dissolved in 20  $\mu$ l TE. The final ligation to produce circles is carried out in a 50  $\mu$ l reaction containing 5  $\mu$ l of v+I linears and 5 units of T4 ligase at 14°C overnight. After 10 min. at 70°C the following day, the reaction mixture is stored at -20°C.

This two-stage procedure results in a molecularly random collection of single-insert plasmid recombinants with minimal contamination from double-insert chimeras (<1%) or free vector (<3%).

Since deviation from randomness can arise from propagation the DNA in the host, *E. coli* host cells deficient in all recombination and restriction functions (A. Grencr, *Strategies 3 (1):5* (1990)) are used to prevent rearrangements, deletions, and loss of clones by restriction. Furthermore, transformed cells are plated directly on antibiotic diffusion plates to avoid the usual broth recovery phase which allows multiplication and selection of the most rapidly growing cells.

Plating is carried out as follows. A 100  $\mu$ l aliquot of Epicurian Coli SURE II Supercompetent Cells (Stratagene 200152) is thawed on ice and transferred to a chilled Falcon 2059 tube on ice. A 1.7  $\mu$ l aliquot of 1.42 M beta-mercaptoethanol is added to the aliquot of cells to a final concentration of 25 mM. Cells are incubated on ice for 10 min. A 1  $\mu$ l aliquot of the final ligation is added to the cells and incubated on ice for 30 min. The cells are heat pulsed for 30 sec. at 42°C and

placed back on ice for 2 min. The outgrowth period in liquid culture is eliminated from this protocol in order to minimize the preferential growth of any given transformed cell. Instead the transformation mixture is plated directly on a nutrient rich SOB plate containing a 5 ml bottom layer of SOB agar (5% SOB agar: 20 g tryptone, 5 g yeast extract, 0.5 g NaCl, 1.5% Difco Agar per liter of media). The 5 ml bottom layer is supplemented with 0.4 ml of 50 mg/ml ampicillin per 100 ml SOB agar. The 15 ml top layer of SOB agar is supplemented with 1 ml X-Gal (2%), 1 ml MgCl<sub>2</sub> (1 M), and 1 ml MgSO<sub>4</sub> /100 ml SOB agar. The 15 ml top layer is poured just prior to plating. Our titer is approximately 100 colonies/10  $\mu$ l aliquot of transformation.<sup>2</sup><sup>4</sup>

All colonies are picked for template preparation regardless of size. Thus, only clones lost due to "poison" DNA or deleterious gene products are deleted from the library, resulting in a slight increase in gap number over that expected.

### 3. Random DNA Sequencing

High quality double stranded DNA plasmid templates are prepared using a "boiling bead" method developed in collaboration with Advanced Genetic Technology Corp. (Gaithersburg, MD) (Adams *et al.*, *Science* 252:1651 (1991); Adams *et al.*, *Nature* 355:632 (1992)). Plasmid preparation is performed in a 96-well format for all stages of DNA preparation from bacterial growth through final DNA purification. Template concentration is determined using Hoechst Dye and a Millipore Cytofluor. DNA concentrations are not adjusted, but low-yielding templates are identified where possible and not sequenced.

Templates are also prepared from two *Streptococcus pneumoniae* lambda genomic libraries. An amplified library is constructed in the vector Lambda GEM-12 (Promega) and an unamplified library is constructed in Lambda DASH II (Stratagene). In particular, for the unamplified lambda library, *Streptococcus pneumoniae* DNA (> 100 kb) is partially digested in a reaction mixture (200  $\mu$ l) containing 50  $\mu$ g DNA, 1X Sau3AI buffer, 20 units Sau3AI for 6 min. at 23°C. The digested DNA was phenol-extracted and electrophoresed on a 0.5% low melting agarose gel at 2V/cm for 7 hours. Fragments from 15 to 25 kb are excised and recovered in a final volume of 6  $\mu$ l. One  $\mu$ l of fragments is used with 1  $\mu$ l of DASHII vector (Stratagene) in the recommended ligation reaction. One  $\mu$ l of the ligation mixture is used per packaging reaction following the recommended protocol with the Gigapack II XL Packaging Extract (Stratagene, #227711). Phage

are plated directly without amplification from the packaging mixture (after dilution with 500  $\mu$ l of recommended SM buffer and chloroform treatment). Yield is about  $2.5 \times 10^3$  pfu/ $\mu$ l. The amplified library is prepared essentially as above except the lambda GEM-12 vector is used. After packaging, about  $3.5 \times 10^4$  pfu are plated on the restrictive NM539 host. The lysate is harvested in 2 ml of SM buffer and stored frozen in 7% dimethylsulfoxide. The phage titer is approximately  $1 \times 10^9$  pfu/ml.

Liquid lysates (100  $\mu$ l) are prepared from randomly selected plaques (from the unamplified library) and template is prepared by long-range PCR using T7 and T3 vector-specific primers.

Sequencing reactions are carried out on plasmid and/or PCR templates using the AB Catalyst LabStation with Applied Biosystems PRISM Ready Reaction Dye Primer Cycle Sequencing Kits for the M13 forward (M13-21) and the M13 reverse (M13RP1) primers (Adams *et al.*, *Nature* 368:474 (1994)). Dye terminator sequencing reactions are carried out on the lambda templates on a Perkin-Elmer 9600 Thermocycler using the Applied Biosystems Ready Reaction Dye Terminator Cycle Sequencing kits. T7 and SP6 primers are used to sequence the ends of the inserts from the Lambda GEM-12 library and T7 and T3 primers are used to sequence the ends of the inserts from the Lambda DASH II library. Sequencing reactions are performed by eight individuals using an average of fourteen AB 373 DNA Sequencers per day. All sequencing reactions are analyzed using the Stretch modification of the AB 373, primarily using a 34 cm well-to-read distance. The overall sequencing success rate very approximately is about 85% for M13-21 and M13RP1 sequences and 65% for dye-terminator reactions. The average usable read length is 485 bp for M13-21 sequences, 445bp for M13RP1 sequences, and 375 bp for dye-terminator reactions.

Richards *et al.*, Chapter 28 in AUTOMATED DNA SEQUENCING AND ANALYSIS, M. D. Adams, C. Fields, J. C. Venter, Eds., Academic Press, London, (1994) described the value of using sequence from both ends of sequencing templates to facilitate ordering of contigs in shotgun assembly projects of lambda and cosmid clones. We balance the desirability of both-end sequencing (including the reduced cost of lower total number of templates) against shorter read-lengths for sequencing reactions performed with the M13RP1 (reverse) primer compared to the M13-21 (forward) primer. Approximately one-half of the templates are sequenced from both ends. Random reverse sequencing reactions are

done based on successful forward sequencing reactions. Some M13RP1 sequences are obtained in a semi-directed fashion: M13-21: sequences pointing outward at the ends of contigs are chosen for M13RP1 sequencing in an effort to specifically order contigs.

5

#### 4. Protocol for Automated Cycle Sequencing

The sequencing is carried out using ABI Catalyst robots and AB 373 Automated DNA Sequencers. The Catalyst robot is a publicly available sophisticated pipetting and temperature control robot which has been developed specifically for DNA sequencing reactions. The Catalyst combines pre-aliquoted templates and reaction mixes consisting of deoxy- and dideoxynucleotides, the thermostable Taq DNA polymerase, fluorescently-labelled sequencing primers, and reaction buffer. Reaction mixes and templates are combined in the wells of an aluminum 96-well thermocycling plate. Thirty consecutive cycles of linear amplification (i.e., one primer synthesis) steps are performed including denaturation, annealing of primer and template, and extension; i.e., DNA synthesis. A heated lid with rubber gaskets on the thermocycling plate prevents evaporation without the need for an oil overlay.

Two sequencing protocols are used: one for dye-labelled primers and a second for dye-labelled dideoxy chain terminators. The shotgun sequencing involves use of four dye-labelled sequencing primers, one for each of the four terminator nucleotide. Each dye-primer is labelled with a different fluorescent dye, permitting the four individual reactions to be combined into one lane of the 373 DNA Sequencer for electrophoresis, detection, and base-calling. ABI currently supplies pre-mixed reaction mixes in bulk packages containing all the necessary non-template reagents for sequencing. Sequencing can be done with both plasmid and PCR-generated templates with both dye-primers and dye-terminators with approximately equal fidelity, although plasmid templates generally give longer usable sequences.

Thirty-two reactions are loaded per AB373 Sequencer each day, for a total of 960 samples. Electrophoresis is run overnight following the manufacturer's protocols, and the data is collected for twelve hours. Following electrophoresis and fluorescence detection, the ABI 373 performs automatic lane tracking and base-calling. The lane-tracking is confirmed visually. Each sequence electropherogram (or fluorescence lane trace) is inspected visually and assessed for quality. Trailing

35

sequences of low quality are removed and the sequence itself is loaded via software to a Sybase database (archived daily to 8mm tape). Leading vector polylinker sequence is removed automatically by a software program. Average edited lengths of sequences from the standard ABI 373 are around 400 bp and depend mostly on the quality of the template used for the sequencing reaction. ABI 373 Sequencers converted to Stretch Lincrs provide a longer electrophoresis path prior to fluorescence detection and increase the average number of usable bases to 500-600 bp.

## INFORMATICS

### 1. Data Management

A number of information management systems for a large-scale sequencing lab have been developed. (For review see, for instance, Kerlavage *et al.*, *Proceedings of the Twenty-Sixth Annual Hawaii International Conference on System Sciences*, IEEE Computer Society Press, Washington D. C., 585 (1993)) The system used to collect and assemble the sequence data was developed using the Sybase relational database management system and was designed to automate data flow wherever possible and to reduce user error. The database stores and correlates all information collected during the entire operation from template preparation to final analysis of the genome. Because the raw output of the ABI 373 Sequencers was based on a Macintosh platform and the data management system chosen was based on a Unix platform, it was necessary to design and implement a variety of multi-user, client-server applications which allow the raw data as well as analysis results to flow seamlessly into the database with a minimum of user effort.

### 2. Assembly

An assembly engine (TIGR Assembler) developed for the rapid and accurate assembly of thousands of sequence fragments was employed to generate contigs. The TIGR assembler simultaneously clusters and assembles fragments of the genome. In order to obtain the speed necessary to assemble more than  $10^4$  fragments, the algorithm builds a hash table of 12 bp oligonucleotide subsequences to generate a list of potential sequence fragment overlaps. The number of potential overlaps for each fragment determines which fragments are likely to fall into repetitive elements. Beginning with a single seed sequence fragment, TIGR Assembler extends the current contig by attempting to add the best matching

fragment based on oligonucleotide content. The contig and candidate fragment are aligned using a modified version of the Smith-Waterman algorithm which provides for optimal gapped alignments (Waterman, M. S., *Methods in Enzymology* 164:765 (1988)). The contig is extended by the fragment only if strict criteria for the quality of the match are met. The match criteria include the minimum length of overlap, the maximum length of an unmatched end, and the minimum percentage match. These criteria are automatically lowered by the algorithm in regions of minimal coverage and raised in regions with a possible repetitive element. The number of potential overlaps for each fragment determines which fragments are likely to fall into repetitive elements. Fragments representing the boundaries of repetitive elements and potentially chimeric fragments are often rejected based on partial mismatches at the ends of alignments and excluded from the current contig. TIGR Assembler is designed to take advantage of clone size information coupled with sequencing from both ends of each template. It enforces the constraint that sequence fragments from two ends of the same template point toward one another in the contig and are located within a certain range of base pairs (definable for each clone based on the known clone size range for a given library).

The process resulted in 391 contigs as represented by SEQ ID NOs:1-391.

### 3. Identifying Genes

The predicted coding regions of the *Streptococcus pneumoniae* genome were initially defined with the program GeneMark, which finds ORFs using a probabilistic classification technique. The predicted coding region sequences were used in searches against a database of all nucleotide sequences from GenBank (October, 1997), using the BLASTN search method to identify overlaps of 50 or more nucleotides with at least a 95% identity. Those ORFs with nucleotide sequence matches are shown in Table 1. The ORFs without such matches were translated to protein sequences and compared to a non-redundant database of known proteins generated by combining the Swiss-prot, PIR and GenPept databases. ORFs that matched a database protein with BLASTP probability less than or equal to 0.01 are shown in Table 2. The table also lists assigned functions based on the closest match in the databases. ORFs that did not match protein or nucleotide sequences in the databases at these levels are shown in Table 3.

## ILLUSTRATIVE APPLICATIONS

### 1. Production of an Antibody to a *Streptococcus pneumoniae* Protein

Substantially pure protein or polypeptide is isolated from the transfected or transformed cells using any one of the methods known in the art. The protein can also be produced in a recombinant prokaryotic expression system, such as *E. coli*, or can be chemically synthesized. Concentration of protein in the final preparation is adjusted, for example, by concentration on an Amicon filter device, to the level of a few micrograms/ml. Monoclonal or polyclonal antibody to the protein can then be prepared as follows.

### 2. Monoclonal Antibody Production by Hybridoma Fusion

Monoclonal antibody to epitopes of any of the peptides identified and isolated as described can be prepared from murine hybridomas according to the classical method of Kohler, G. and Milstein, C., *Nature* 256:495 (1975) or modifications of the methods thereof. Briefly, a mouse is repetitively inoculated with a few micrograms of the selected protein over a period of a few weeks. The mouse is then sacrificed, and the antibody producing cells of the spleen isolated. The spleen cells are fused by means of polyethylene glycol with mouse myeloma cells, and the excess unfused cells destroyed by growth of the system on selective media comprising aminopterin (HAT media). The successfully fused cells are diluted and aliquots of the dilution placed in wells of a microtiter plate where growth of the culture is continued. Antibody-producing clones are identified by detection of antibody in the supernatant fluid of the wells by immunoassay procedures, such as ELISA, as originally described by Engvall, E., *Meih. Enzymol.* 70:419 (1980), and modified methods thereof. Selected positive clones can be expanded and their monoclonal antibody product harvested for use. Detailed procedures for monoclonal antibody production are described in Davis, L. *et al.*, *Basic Methods in Molecular Biology*, Elsevier, New York. Section 21-2 (1989).

### 3. Polyclonal Antibody Production by Immunization

Polyclonal antiserum containing antibodies to heterogeneous epitopes of a single protein can be prepared by immunizing suitable animals with the expressed protein described above, which can be unmodified or modified to enhance immunogenicity. Effective polyclonal antibody production is affected by many factors related both to the antigen and the host species. For example, small molecules tend to be less immunogenic than others and may require the use of carriers and adjuvant. Also, host animals vary in response to site of inoculations and dose, with both inadequate or excessive doses of antigen resulting in low titer antisera. Small doses (ng level) of antigen administered at multiple intradermal sites appears to be most reliable. An effective immunization protocol for rabbits can be found in Vaitukaitis, J. *et al.*, *J. Clin. Endocrinol. Metab.* 33:988-991 (1971).

Booster injections can be given at regular intervals, and antiscrum harvested when antibody titer thereof, as determined semi-quantitatively, for example, by double immunodiffusion in agar against known concentrations of the antigen, begins to fall. See, for example, Ouchterlony, O. *et al.*, Chap. 19 in: *Handbook of Experimental Immunology*, Wier, D., ed, Blackwell (1973). Plateau concentration of antibody is usually in the range of 0.1 to 0.2 mg/ml of serum (about 12M). Affinity of the antisera for the antigen is determined by preparing competitive binding curves, as described, for example, by Fisher, D., Chap. 42 in: *Manual of Clinical Immunology*, second edition, Rose and Friedman, eds., Amer. Soc. For Microbiology, Washington, D. C. (1980)

Antibody preparations prepared according to either protocol are useful in quantitative immunoassays which determine concentrations of antigen-bearing substances in biological samples; they are also used semi- quantitatively or qualitatively to identify the presence of antigen in a biological sample. In addition, antibodies are useful in various animal models of pneumococcal disease as a means of evaluating the protein used to make the antibody as a potential vaccine target or as a means of evaluating the antibody as a potential immunotherapeutic or immunoprophylactic reagent.



#### 4. Preparation of PCR Primers and Amplification of DNA

Various fragments of the *Streptococcus pneumoniae* genome, such as those of Tables 1-3 and SEQ ID NOS:1-391 can be used, in accordance with the present invention, to prepare PCR primers for a variety of uses. The PCR primers are preferably at least 15 bases, and more preferably at least 18 bases in length. When selecting a primer sequence, it is preferred that the primer pairs have approximately the same G/C ratio, so that melting temperatures are approximately the same. The PCR primers and amplified DNA of this Example find use in the Examples that follow.

#### 5. Gene expression from DNA Sequences Corresponding to ORFs

A fragment of the *Streptococcus pneumoniae* genome provided in Tables 1-3 is introduced into an expression vector using conventional technology. Techniques to transfer cloned sequences into expression vectors that direct protein translation in mammalian, yeast, insect or bacterial expression systems are well known in the art. Commercially available vectors and expression systems are available from a variety of suppliers including Stratagene (La Jolla, California), Promega (Madison, Wisconsin), and Invitrogen (San Diego, California). If desired, to enhance expression and facilitate proper protein folding, the codon context and codon pairing of the sequence may be optimized for the particular expression organism, as explained by Hatfield *et al.*, U. S. Patent No. 5,082,767, incorporated herein by this reference.

The following is provided as one exemplary method to generate polypeptide(s) from cloned ORFs of the *Streptococcus pneumoniae* genome fragment. Bacterial ORFs generally lack a poly A addition signal. The addition signal sequence can be added to the construct by, for example, splicing out the poly A addition sequence from pSG5 (Stratagene) using BglI and SalI restriction endonuclease enzymes and incorporating it into the mammalian expression vector pXT1 (Stratagene) for use in eukaryotic expression systems. pXT1 contains the LTRs and a portion of the gag gene of Moloney Murine Leukemia Virus. The positions of the LTRs in the construct allow efficient stable transfection. The vector includes the Herpes Simplex thymidine kinase promoter and the selectable neomycin gene. The *Streptococcus pneumoniae* DNA is obtained by PCR from the bacterial vector using oligonucleotide primers complementary to the *Streptococcus pneumoniae* DNA and containing restriction endonuclease sequences for PstI incorporated into the 5' primer and BglII at the 5' end of the corresponding *Streptococcus pneumoniae* DNA 3' primer, taking care to ensure that the *Streptococcus pneumoniae* DNA is positioned such that its followed with the poly A addition sequence. The purified fragment obtained from the resulting PCR reaction is digested with PstI, blunt ended with an exonuclease, digested with BglII, purified and ligated to pXT1, now containing a poly A addition sequence and digested BglII.

The ligated product is transfected into mouse NIH 3T3 cells using Lipofectin (Life Technologies, Inc., Grand Island, New York) under conditions outlined in the product specification. Positive transfectants are selected after growing the transfected cells in 600 ug/ml G418 (Sigma, St. Louis, Missouri). The protein is preferably released into the supernatant. However if the protein has membrane binding domains, the protein may additionally be retained within the cell or expression may be restricted to the cell surface. Since it may be necessary to purify and locate the transfected product, synthetic 15-mer peptides synthesized from the predicted *Streptococcus pneumoniae* DNA sequence are injected into mice to generate antibody to the polypeptide encoded by the *Streptococcus pneumoniae* DNA.

Alternatively and if antibody production is not possible, the *Streptococcus pneumoniae* DNA sequence is additionally incorporated into eukaryotic expression vectors and expressed as, for example, a globin fusion. Antibody to the globin moiety then is used to purify the chimeric protein. Corresponding protease  
5 cleavage sites are engineered between the globin moiety and the polypeptide encoded by the *Streptococcus pneumoniae* DNA so that the latter may be freed from the formed by simple protease digestion. One useful expression vector for generating globin chimerics is pSG5 (Stratagene). This vector encodes a rabbit globin. Intron II of the rabbit globin gene facilitates splicing of the expressed  
10 transcript, and the polyadenylation signal incorporated into the construct increases the level of expression. These techniques are well known to those skilled in the art of molecular biology. Standard methods are published in methods texts such as Davis *et al.*, cited elsewhere herein, and many of the methods are available from the technical assistance representatives from Stratagene, Life Technologies, Inc., or  
15 Promega. Polypeptides of the invention also may be produced using *in vitro* translation systems such as *in vitro* Express™ Translation Kit (Stratagene).

While the present invention has been described in some detail for purposes of clarity and understanding, one skilled in the art will appreciate that various changes in form and detail can be made without departing from the true scope of  
20 the invention.

All patents, patent applications and publications referred to above are hereby incorporated by reference.

TABLE 1

S. pneumoniae - Coding regions containing known sequences

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	percent Ident	RSP nt length	ORF nt length
1	1	437	1003	gb U41735	Streptococcus pneumoniae 5S2 rRNA gene, complete cds	92	200	567
2	5	6169	5720	gb U40471	Streptococcus pneumoniae 5S2 rRNA gene, complete cds	96	450	450
2	6	6592	6167	emb 283335 SP28	Streptococcus pneumoniae 5S2 rRNA gene, complete cds	98	426	426
3	11	9770	9147	emb 283335 SP28	S. pneumoniae desB, cap1A, B, C, D, E, F, G, H, I, J, K genes, dTDP-thiamase biosynthesis genes and a1a gene	94	624	624
3	12	10489	9671	emb 283335 SP28	S. pneumoniae desB, cap1A, B, C, D, E, F, G, H, I, J, K genes, dTDP-thiamase biosynthesis genes and a1a gene	91	819	819
3	13	11546	12019	gb U43526	Streptococcus pneumoniae neuraminidase B (nanB) gene, complete cds, and neuraminidase (nanA) gene, partial cds	99	474	474
3	14	12017	13375	gb U43526	Streptococcus pneumoniae neuraminidase B (nanB) gene, complete cds, and neuraminidase (nanA) gene, partial cds	99	1359	1359
3	15	13421	14338	gb U43526	Streptococcus pneumoniae neuraminidase B (nanB) gene, complete cds, and neuraminidase (nanA) gene, partial cds	99	918	918
3	16	14329	15171	gb U43526	Streptococcus pneumoniae neuraminidase B (nanB) gene, complete cds, and neuraminidase (nanA) gene, partial cds	99	843	843
3	17	15132	17282	gb U43526	Streptococcus pneumoniae neuraminidase B (nanB) gene, complete cds, and neuraminidase (nanA) gene, partial cds	99	2151	2151
3	18	17267	18397	gb U43526	Streptococcus pneumoniae neuraminidase B (nanB) gene, complete cds, and neuraminidase (nanA) gene, partial cds	99	1069	1131
4	1	46	1186	emb Y11463 SP20	Streptococcus pneumoniae dng, rpoD, cpdA genes and ORF3 and ORF5	99	1143	1143
4	2	1198	2529	emb Y11463 SP20	Streptococcus pneumoniae dng, rpoD, cpdA genes and ORF3 and ORF5	99	876	1332
5	7	11297	11473	gb U41735	Streptococcus pneumoniae peptide methionine sulfoxide reductase (metA) and methionine sulfoxide reductase (metB) genes, complete cds	82	175	177
6	7	7125	7364	emb 277726 SP15	S. pneumoniae DNA for insertion sequence IS1318 (1372 bp)	93	238	240
6	8	7322	7570	emb 277726 SP15	S. pneumoniae DNA for insertion sequence IS1318 (1372 bp)	95	160	249
6	9	7533	7985	emb 277726 SP15	S. pneumoniae DNA for insertion sequence IS1318 (1372 bp)	99	453	453
6	23	20197	19733	emb 283335 SP28	S. pneumoniae desB, cap1A, B, C, D, E, F, G, H, I, J, K genes, dTDP-thiamase biosynthesis genes and a1a gene	96	465	465
7	10	8305	7682	emb 283335 SP28	S. pneumoniae desB, cap1A, B, C, D, E, F, G, H, I, J, K genes, dTDP-thiamase biosynthesis genes and a1a gene	95	624	624

TABLE 1

S. pneumoniae - Coding regions containing known sequences

Coding region	ORF no.	Start no.	Stop no.	match region	match gene name	percent ident	RSP nt length	ORF nt length
7	11	9024	8206	emb 283335 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-ribose biosynthesis genes and allia gene	95	819	819
10	13	9104	8078	gb 129323	Streptococcus pneumoniae methyl transferase (acr) gene cluster, complete cds	93	513	1227
11	2	346	319	emb 279691 SOOR	S.pneumoniae yor(A,B,C,D,E), ftsL, pbpX and regK genes	99	316	372
11	3	892	1980	emb 279691 SOOR	S.pneumoniae yor(A,B,C,D,E), ftsL, pbpX and regK genes	99	1089	1089
11	5	3040	3477	emb 279691 SOOR	S.pneumoniae yor(A,B,C,D,E), ftsL, pbpX and regK genes	99	259	438
11	6	3480	3247	emb 279691 SOOR	S.pneumoniae yor(A,B,C,D,E), ftsL, pbpX and regK genes	99	234	234
11	7	3601	4537	emb 279691 SOOR	S.pneumoniae yor(A,B,C,D,E), ftsL, pbpX and regK genes	98	957	957
11	8	4506	4886	emb 279691 SOOR	S.pneumoniae yor(A,B,C,D,E), ftsL, pbpX and regK genes	99	381	381
11	9	4884	7142	emb X16367 SPPB	Streptococcus pneumoniae pbpX gene for penicillin binding protein 2X	99	2259	2259
11	10	7132	8134	emb X16367 SPPB	Streptococcus pneumoniae pbpX gene for penicillin binding protein 2X	98	70	993
13	1	53	1126	gb H31296	S.pneumoniae recP gene, complete cds	99	437	1074
14	3	1837	2148	emb 283335 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-ribose biosynthesis genes and allia gene	87	96	312
14	4	2318	2108	gb H36180	Streptococcus pneumoniae Transposase, (comA and comB) and SAICAR synthetase partial cds	98	411	411
15	9	8942	8511	gb U09239	Streptococcus pneumoniae type 19F capsular polysaccharide biosynthesis operon, cpsA(cpsA2)cpsH(cpsH2) genes, complete cds, and allia gene,	89	340	432
17	7	3910	3458	emb 277226 SP16	S.pneumoniae DNA for insertion sequence IS1318 (1372 bp)	98	453	453
17	8	4304	3873	emb 277227 SP16	S.pneumoniae DNA for insertion sequence IS1318 (823 bp)	96	382	432
19	1	41	529	emb X94509 SP10	S.pneumoniae iga gene	75	560	489
19	2	554	157	gb U07722	Streptococcus pneumoniae attachment site (actB), DNA sequence	99	167	204
19	3	946	1877	gb U07721	Streptococcus pneumoniae attachment site (actB), DNA sequence	94	100	882
20	1	937	187	gb U33315	Streptococcus pneumoniae orfL gene, partial cds, competence stimulating protein (comC), histidine protein kinase (comD) and response regulator (comE) genes	99	756	756
20	2	2271	931	gb U33315	Streptococcus pneumoniae orfL gene, partial cds, competence stimulating protein (comC), histidine protein kinase (comD) and response regulator (comE) genes	98	1341	1341



TABLE I

S. pneumoniae - Coding regions containing known sequences

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match gene name	Percent Ident	HSP nt length	ORF nt length
26	8	14498	14854	emb 283335 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and allia gene	99	338	357
26	9	14763	14924	emb 283335 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and allia gene	100	94	162
26	10	14922	15173	gb U04047	Streptococcus pneumoniae SS2 dextran glucosidase gene and Insertion sequence 181202 transposase gene, complete cds	97	242	252
28	1	80	505	emb 283335 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and allia gene	99	426	426
28	2	503	952	gb U04047	Streptococcus pneumoniae SS2 dextran glucosidase gene and Insertion sequence 181202 transposase gene, complete cds	97	450	450
28	3	780	1238	gb U04047	Streptococcus pneumoniae SS2 dextran glucosidase gene and Insertion sequence 181202 transposase gene, complete cds	96	181	519
34	1	207	1323	gb U06611	Streptococcus pneumoniae maltose/maltodextrin uptake (malK) and two maltodextrin permease (malC and malD) genes, complete cds	99	1317	1317
34	2	1477	2367	gb U06611	Streptococcus pneumoniae maltose/maltodextrin uptake (malK) and two maltodextrin permease (malC and malD) genes, complete cds	96	795	891
34	3	2593	3420	gb U21855	Streptococcus pneumoniae malK gene, complete cds; malR gene, complete cds	96	446	828
34	4	2790	2867	gb U21855	Streptococcus pneumoniae malK gene, complete cds; malR gene, complete cds	96	137	144
34	5	3418	4416	gb U21855	Streptococcus pneumoniae malK gene, complete cds; malR gene, complete cds	96	999	999
34	9	7764	7907	gb U41735	Streptococcus pneumoniae peptidase methionine sulfoxide reductase (metA) and homoserine kinase homolog (thrB) genes, complete cds	93	201	258
34	16	10562	10257	emb X63602 SPPO	S.pneumoniae malK-box			
35	4	1176	1439	emb 283335 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and allia gene	92	238	306
35	5	1458	1961	gb U09239	Streptococcus pneumoniae type 19F capsular polysaccharide biosynthesis operon, (cpsA)(cpsB)(cpsC)(cpsD)(cpsE)(cpsF)(cpsG)(cpsH)(cpsI)(cpsJ)(cpsK)(cpsL)(cpsM)(cpsN)(cpsO)(cpsP)(cpsQ)(cpsR)(cpsS)(cpsT)(cpsU)(cpsV)(cpsW)(cpsX)(cpsY)(cpsZ)(cpsAA)(cpsAB)(cpsAC)(cpsAD)(cpsAE)(cpsAF)(cpsAG)(cpsAH)(cpsAI)(cpsAJ)(cpsAK)(cpsAL)(cpsAM)(cpsAN)(cpsAO)(cpsAP)(cpsAQ)(cpsAR)(cpsAS)(cpsAT)(cpsAU)(cpsAV)(cpsAW)(cpsAX)(cpsAY)(cpsAZ)(cpsBA)(cpsBB)(cpsBC)(cpsBD)(cpsBE)(cpsBF)(cpsBG)(cpsBH)(cpsBI)(cpsBJ)(cpsBK)(cpsBL)(cpsBM)(cpsBN)(cpsBO)(cpsBP)(cpsBQ)(cpsBR)(cpsBS)(cpsBT)(cpsBU)(cpsBV)(cpsBW)(cpsBX)(cpsBY)(cpsBZ)(cpsCA)(cpsCB)(cpsCC)(cpsCD)(cpsCE)(cpsCF)(cpsCG)(cpsCH)(cpsCI)(cpsCJ)(cpsCK)(cpsCL)(cpsCM)(cpsCN)(cpsCO)(cpsCP)(cpsCQ)(cpsCR)(cpsCS)(cpsCT)(cpsCU)(cpsCV)(cpsCW)(cpsCX)(cpsCY)(cpsCZ)(cpsDA)(cpsDB)(cpsDC)(cpsDD)(cpsDE)(cpsDF)(cpsDG)(cpsDH)(cpsDI)(cpsDJ)(cpsDK)(cpsDL)(cpsDM)(cpsDN)(cpsDO)(cpsDP)(cpsDQ)(cpsDR)(cpsDS)(cpsDT)(cpsDU)(cpsDV)(cpsDW)(cpsDX)(cpsDY)(cpsDZ)(cpsEA)(cpsEB)(cpsEC)(cpsED)(cpsEE)(cpsEF)(cpsEG)(cpsEH)(cpsEI)(cpsEJ)(cpsEK)(cpsEL)(cpsEM)(cpsEN)(cpsEO)(cpsEP)(cpsEQ)(cpsER)(cpsES)(cpsET)(cpsEU)(cpsEV)(cpsEW)(cpsEX)(cpsEY)(cpsEZ)(cpsFA)(cpsFB)(cpsFC)(cpsFD)(cpsFE)(cpsFF)(cpsFG)(cpsFH)(cpsFI)(cpsFJ)(cpsFK)(cpsFL)(cpsFM)(cpsFN)(cpsFO)(cpsFP)(cpsFQ)(cpsFR)(cpsFS)(cpsFT)(cpsFU)(cpsFV)(cpsFW)(cpsFX)(cpsFY)(cpsFZ)(cpsGA)(cpsGB)(cpsGC)(cpsGD)(cpsGE)(cpsGF)(cpsGG)(cpsGH)(cpsGI)(cpsGJ)(cpsGK)(cpsGL)(cpsGM)(cpsGN)(cpsGO)(cpsGP)(cpsGQ)(cpsGR)(cpsGS)(cpsGT)(cpsGU)(cpsGV)(cpsGW)(cpsGX)(cpsGY)(cpsGZ)(cpsHA)(cpsHB)(cpsHC)(cpsHD)(cpsHE)(cpsHF)(cpsHG)(cpsHH)(cpsHI)(cpsHJ)(cpsHK)(cpsHL)(cpsHM)(cpsHN)(cpsHO)(cpsHP)(cpsHQ)(cpsHR)(cpsHS)(cpsHT)(cpsHU)(cpsHV)(cpsHW)(cpsHX)(cpsHY)(cpsHZ)(cpsIA)(cpsIB)(cpsIC)(cpsID)(cpsIE)(cpsIF)(cpsIG)(cpsIH)(cpsII)(cpsIJ)(cpsIK)(cpsIL)(cpsIM)(cpsIN)(cpsIO)(cpsIP)(cpsIQ)(cpsIR)(cpsIS)(cpsIT)(cpsIU)(cpsIV)(cpsIW)(cpsIX)(cpsIY)(cpsIZ)(cpsJA)(cpsJB)(cpsJC)(cpsJD)(cpsJE)(cpsJF)(cpsJG)(cpsJH)(cpsJI)(cpsJJ)(cpsJK)(cpsJL)(cpsJM)(cpsJN)(cpsJO)(cpsJP)(cpsJQ)(cpsJR)(cpsJS)(cpsJT)(cpsJU)(cpsJV)(cpsJW)(cpsJX)(cpsJY)(cpsJZ)(cpsKA)(cpsKB)(cpsKC)(cpsKD)(cpsKE)(cpsKF)(cpsKG)(cpsKH)(cpsKI)(cpsKJ)(cpsKL)(cpsKM)(cpsKN)(cpsKO)(cpsKP)(cpsKQ)(cpsKR)(cpsKS)(cpsKT)(cpsKU)(cpsKV)(cpsKW)(cpsKX)(cpsKY)(cpsKZ)(cpsLA)(cpsLB)(cpsLC)(cpsLD)(cpsLE)(cpsLF)(cpsLG)(cpsLH)(cpsLI)(cpsLJ)(cpsLK)(cpsLL)(cpsLM)(cpsLN)(cpsLO)(cpsLP)(cpsLQ)(cpsLR)(cpsLS)(cpsLT)(cpsLU)(cpsLV)(cpsLW)(cpsLX)(cpsLY)(cpsLZ)(cpsMA)(cpsMB)(cpsMC)(cpsMD)(cpsME)(cpsMF)(cpsMG)(cpsMH)(cpsMI)(cpsMJ)(cpsMK)(cpsML)(cpsMN)(cpsMO)(cpsMP)(cpsMQ)(cpsMR)(cpsMS)(cpsMT)(cpsMU)(cpsMV)(cpsMW)(cpsMX)(cpsMY)(cpsMZ)(cpsNA)(cpsNB)(cpsNC)(cpsND)(cpsNE)(cpsNF)(cpsNG)(cpsNH)(cpsNI)(cpsNJ)(cpsNK)(cpsNL)(cpsNM)(cpsNO)(cpsNP)(cpsNQ)(cpsNR)(cpsNS)(cpsNT)(cpsNU)(cpsNV)(cpsNW)(cpsNX)(cpsNY)(cpsNZ)(cpsOA)(cpsOB)(cpsOC)(cpsOD)(cpsOE)(cpsOF)(cpsOG)(cpsOH)(cpsOI)(cpsOJ)(cpsOK)(cpsOL)(cpsOM)(cpsON)(cpsOO)(cpsOP)(cpsOQ)(cpsOR)(cpsOS)(cpsOT)(cpsOU)(cpsOV)(cpsOW)(cpsOX)(cpsOY)(cpsOZ)(cpsPA)(cpsPB)(cpsPC)(cpsPD)(cpsPE)(cpsPF)(cpsPG)(cpsPH)(cpsPI)(cpsPJ)(cpsPK)(cpsPL)(cpsPM)(cpsPN)(cpsPO)(cpsPP)(cpsPQ)(cpsPR)(cpsPS)(cpsPT)(cpsPU)(cpsPV)(cpsPW)(cpsPX)(cpsPY)(cpsPZ)(cpsQA)(cpsQB)(cpsQC)(cpsQD)(cpsQE)(cpsQF)(cpsQG)(cpsQH)(cpsQI)(cpsQJ)(cpsQK)(cpsQL)(cpsQM)(cpsQN)(cpsQO)(cpsQP)(cpsQQ)(cpsQR)(cpsQS)(cpsQT)(cpsQU)(cpsQV)(cpsQW)(cpsQX)(cpsQY)(cpsQZ)(cpsRA)(cpsRB)(cpsRC)(cpsRD)(cpsRE)(cpsRF)(cpsRG)(cpsRH)(cpsRI)(cpsRJ)(cpsRK)(cpsRL)(cpsRM)(cpsRN)(cpsRO)(cpsRP)(cpsRQ)(cpsRR)(cpsRS)(cpsRT)(cpsRU)(cpsRV)(cpsRW)(cpsRX)(cpsRY)(cpsRZ)(cpsSA)(cpsSB)(cpsSC)(cpsSD)(cpsSE)(cpsSF)(cpsSG)(cpsSH)(cpsSI)(cpsSJ)(cpsSK)(cpsSL)(cpsSM)(cpsSN)(cpsSO)(cpsSP)(cpsSQ)(cpsSR)(cpsSS)(cpsST)(cpsSU)(cpsSV)(cpsSW)(cpsSX)(cpsSY)(cpsSZ)(cpsTA)(cpsTB)(cpsTC)(cpsTD)(cpsTE)(cpsTF)(cpsTG)(cpsTH)(cpsTI)(cpsTJ)(cpsTK)(cpsTL)(cpsTM)(cpsTN)(cpsTO)(cpsTP)(cpsTQ)(cpsTR)(cpsTS)(cpsTT)(cpsTU)(cpsTV)(cpsTW)(cpsTX)(cpsTY)(cpsTZ)(cpsUA)(cpsUB)(cpsUC)(cpsUD)(cpsUE)(cpsUF)(cpsUG)(cpsUH)(cpsUI)(cpsUJ)(cpsUK)(cpsUL)(cpsUM)(cpsUN)(cpsUO)(cpsUP)(cpsUQ)(cpsUR)(cpsUS)(cpsUT)(cpsUU)(cpsUV)(cpsUW)(cpsUX)(cpsUY)(cpsUZ)(cpsVA)(cpsVB)(cpsVC)(cpsVD)(cpsVE)(cpsVF)(cpsVG)(cpsVH)(cpsVI)(cpsVJ)(cpsVK)(cpsVL)(cpsVM)(cpsVN)(cpsVO)(cpsVP)(cpsVQ)(cpsVR)(cpsVS)(cpsVT)(cpsVU)(cpsVV)(cpsVW)(cpsVX)(cpsVY)(cpsVZ)(cpsWA)(cpsWB)(cpsWC)(cpsWD)(cpsWE)(cpsWF)(cpsWG)(cpsWH)(cpsWI)(cpsWJ)(cpsWK)(cpsWL)(cpsWM)(cpsWN)(cpsWO)(cpsWP)(cpsWQ)(cpsWR)(cpsWS)(cpsWT)(cpsWU)(cpsWV)(cpsWW)(cpsWX)(cpsWY)(cpsWZ)(cpsXA)(cpsXB)(cpsXC)(cpsXD)(cpsXE)(cpsXF)(cpsXG)(cpsXH)(cpsXI)(cpsXJ)(cpsXK)(cpsXL)(cpsXM)(cpsXN)(cpsXO)(cpsXP)(cpsXQ)(cpsXR)(cpsXS)(cpsXT)(cpsXU)(cpsXV)(cpsXW)(cpsXX)(cpsXY)(cpsXZ)(cpsYA)(cpsYB)(cpsYC)(cpsYD)(cpsYE)(cpsYF)(cpsYG)(cpsYH)(cpsYI)(cpsYJ)(cpsYK)(cpsYL)(cpsYM)(cpsYN)(cpsYO)(cpsYP)(cpsYQ)(cpsYR)(cpsYS)(cpsYT)(cpsYU)(cpsYV)(cpsYW)(cpsYX)(cpsYY)(cpsYZ)(cpsZA)(cpsZB)(cpsZC)(cpsZD)(cpsZE)(cpsZF)(cpsZG)(cpsZH)(cpsZI)(cpsZJ)(cpsZK)(cpsZL)(cpsZM)(cpsZN)(cpsZO)(cpsZP)(cpsZQ)(cpsZR)(cpsZS)(cpsZT)(cpsZU)(cpsZV)(cpsZW)(cpsZX)(cpsZY)(cpsZZ)	98	264	504
35	17	16172	15477	emb X63707 SPCP	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and allia gene	97	696	696
35	18	16961	16120	emb 283335 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and allia gene	86	792	792
35	19	17620	16671	gb U09239	Streptococcus pneumoniae type 19F capsular polysaccharide biosynthesis operon, (cpsA)(cpsB)(cpsC)(cpsD)(cpsE)(cpsF)(cpsG)(cpsH)(cpsI)(cpsJ)(cpsK)(cpsL)(cpsM)(cpsN)(cpsO)(cpsP)(cpsQ)(cpsR)(cpsS)(cpsT)(cpsU)(cpsV)(cpsW)(cpsX)(cpsY)(cpsZ)(cpsAA)(cpsAB)(cpsAC)(cpsAD)(cpsAE)(cpsAF)(cpsAG)(cpsAH)(cpsAI)(cpsAJ)(cpsAK)(cpsAL)(cpsAM)(cpsAN)(cpsAO)(cpsAP)(cpsAQ)(cpsAR)(cpsAS)(cpsAT)(cpsAU)(cpsAV)(cpsAW)(cpsAX)(cpsAY)(cpsAZ)(cpsBA)(cpsBB)(cpsBC)(cpsBD)(cpsBE)(cpsBF)(cpsBG)(cpsBH)(cpsBI)(cpsBJ)(cpsBK)(cpsBL)(cpsBM)(cpsBN)(cpsBO)(cpsBP)(cpsBQ)(cpsBR)(cpsBS)(cpsBT)(cpsBU)(cpsBV)(cpsBW)(cpsBX)(cpsBY)(cpsBZ)(cpsCA)(cpsCB)(cpsCC)(cpsCD)(cpsCE)(cpsCF)(cpsCG)(cpsCH)(cpsCI)(cpsCJ)(cpsCK)(cpsCL)(cpsCM)(cpsCN)(cpsCO)(cpsCP)(cpsCQ)(cpsCR)(cpsCS)(cpsCT)(cpsCU)(cpsCV)(cpsCW)(cpsCX)(cpsCY)(cpsCZ)(cpsDA)(cpsDB)(cpsDC)(cpsDD)(cpsDE)(cpsDF)(cpsDG)(cpsDH)(cpsDI)(cpsDJ)(cpsDK)(cpsDL)(cpsDM)(cpsDN)(cpsDO)(cpsDP)(cpsDQ)(cpsDR)(cpsDS)(cpsDT)(cpsDU)(cpsDV)(cpsDW)(cpsDX)(cpsDY)(cpsDZ)(cpsEA)(cpsEB)(cpsEC)(cpsED)(cpsEE)(cpsEF)(cpsEG)(cpsEH)(cpsEI)(cpsEJ)(cpsEK)(cpsEL)(cpsEM)(cpsEN)(cpsEO)(cpsEP)(cpsEQ)(cpsER)(cpsES)(cpsET)(cpsEU)(cpsEV)(cpsEW)(cpsEX)(cpsEY)(cpsEZ)(cpsFA)(cpsFB)(cpsFC)(cpsFD)(cpsFE)(cpsFF)(cpsFG)(cpsFH)(cpsFI)(cpsFJ)(cpsFK)(cpsFL)(cpsFM)(cpsFN)(cpsFO)(cpsFP)(cpsFQ)(cpsFR)(cpsFS)(cpsFT)(cpsFU)(cpsFV)(cpsFW)(cpsFX)(cpsFY)(cpsFZ)(cpsGA)(cpsGB)(cpsGC)(cpsGD)(cpsGE)(cpsGF)(cpsGG)(cpsGH)(cpsGI)(cpsGJ)(cpsGK)(cpsGL)(cpsGM)(cpsGN)(cpsGO)(cpsGP)(cpsGQ)(cpsGR)(cpsGS)(cpsGT)(cpsGU)(cpsGV)(cpsGW)(cpsGX)(cpsGY)(cpsGZ)(cpsHA)(cpsHB)(cpsHC)(cpsHD)(cpsHE)(cpsHF)(cpsHG)(cpsHH)(cpsHI)(cpsHJ)(cpsHK)(cpsHL)(cpsHM)(cpsHN)(cpsHO)(cpsHP)(cpsHQ)(cpsHR)(cpsHS)(cpsHT)(cpsHU)(cpsHV)(cpsHW)(cpsHX)(cpsHY)(cpsHZ)(cpsIA)(cpsIB)(cpsIC)(cpsID)(cpsIE)(cpsIF)(cpsIG)(cpsIH)(cpsII)(cpsIJ)(cpsIK)(cpsIL)(cpsIM)(cpsIN)(cpsIO)(cpsIP)(cpsIQ)(cpsIR)(cpsIS)(cpsIT)(cpsIU)(cpsIV)(cpsIW)(cpsIX)(cpsIY)(cpsIZ)(cpsJA)(cpsJB)(cpsJC)(cpsJD)(cpsJE)(cpsJF)(cpsJG)(cpsJH)(cpsJI)(cpsJJ)(cpsJK)(cpsJL)(cpsJM)(cpsJN)(cpsJO)(cpsJP)(cpsJQ)(cpsJR)(cpsJS)(cpsJT)(cpsJU)(cpsJV)(cpsJW)(cpsJX)(cpsJY)(cpsJZ)(cpsKA)(cpsKB)(cpsKC)(cpsKD)(cpsKE)(cpsKF)(cpsKG)(cpsKH)(cpsKI)(cpsKJ)(cpsKL)(cpsKM)(cpsKN)(cpsKO)(cpsKP)(cpsKQ)(cpsKR)(cpsKS)(cpsKT)(cpsKU)(cpsKV)(cpsKW)(cpsKX)(cpsKY)(cpsKZ)(cpsLA)(cpsLB)(cpsLC)(cpsLD)(cpsLE)(cpsLF)(cpsLG)(cpsLH)(cpsLI)(cpsLJ)(cpsLK)(cpsLM)(cpsLN)(cpsLO)(cpsLP)(cpsLQ)(cpsLR)(cpsLS)(cpsLT)(cpsLU)(cpsLV)(cpsLW)(cpsLX)(cpsLY)(cpsLZ)(cpsMA)(cpsMB)(cpsMC)(cpsMD)(cpsME)(cpsMF)(cpsMG)(cpsMH)(cpsMI)(cpsMJ)(cpsMK)(cpsML)(cpsMN)(cpsMO)(cpsMP)(cpsMQ)(cpsMR)(cpsMS)(cpsMT)(cpsMU)(cpsMV)(cpsMW)(cpsMX)(cpsMY)(cpsMZ)(cpsNA)(cpsNB)(cpsNC)(cpsND)(cpsNE)(cpsNF)(cpsNG)(cpsNH)(cpsNI)(cpsNJ)(cpsNK)(cpsNL)(cpsNM)(cpsNO)(cpsNP)(cpsNQ)(cpsNR)(cpsNS)(cpsNT)(cpsNU)(cpsNV)(cpsNW)(cpsNX)(cpsNY)(cpsNZ)(cpsOA)(cpsOB)(cpsOC)(cpsOD)(cpsOE)(cpsOF)(cpsOG)(cpsOH)(cpsOI)(cpsOJ)(cpsOK)(cpsOL)(cpsOM)(cpsON)(cpsOO)(cpsOP)(cpsOQ)(cpsOR)(cpsOS)(cpsOT)(cpsOU)(cpsOV)(cpsOW)(cpsOX)(cpsOY)(cpsOZ)(cpsPA)(cpsPB)(cpsPC)(cpsPD)(cpsPE)(cpsPF)(cpsPG)(cpsPH)(cpsPI)(cpsPJ)(cpsPK)(cpsPL)(cpsPM)(cpsPN)(cpsPO)(cpsPP)(cpsPQ)(cpsPR)(cpsPS)(cpsPT)(cpsPU)(cpsPV)(cpsPW)(cpsPX)(cpsPY)(cpsPZ)(cpsQA)(cpsQB)(cpsQC)(cpsQD)(cpsQE)(cpsQF)(cpsQG)(cpsQH)(cpsQI)(cpsQJ)(cpsQK)(cpsQL)(cpsQM)(cpsQN)(cpsQO)(cpsQP)(cpsQQ)(cpsQR)(cpsQS)(cpsQT)(cpsQU)(cpsQV)(cpsQW)(cpsQX)(cpsQY)(cpsQZ)(cpsRA)(cpsRB)(cpsRC)(cpsRD)(cpsRE)(cpsRF)(cpsRG)(cpsRH)(cpsRI)(cpsRJ)(cpsRK)(cpsRL)(cpsRM)(cpsRN)(cpsRO)(cpsRP)(cpsRQ)(cpsRR)(cpsRS)(cpsRT)(cpsRU)(cpsRV)(cpsRW)(cpsRX)(cpsRY)(cpsRZ)(cpsSA)(cpsSB)(cpsSC)(cpsSD)(cpsSE)(cpsSF)(cpsSG)(cpsSH)(cpsSI)(cpsSJ)(cpsSK)(cpsSL)(cpsSM)(cpsSN)(cpsSO)(cpsSP)(cpsSQ)(cpsSR)(cpsSS)(cpsST)(cpsSU)(cpsSV)(cpsSW)(cpsSX)(cpsSY)(cpsSZ)(cpsTA)(cpsTB)(cpsTC)(cpsTD)(cpsTE)(cpsTF)(cpsTG)(cpsTH)(cpsTI)(cpsTJ)(cpsTK)(cpsTL)(cpsTM)(cpsTN)(cpsTO)(cpsTP)(cpsTQ)(cpsTR)(cpsTS)(cpsTT)(cpsTU)(cpsTV)(cpsTW)(cpsTX)(cpsTY)(cpsTZ)(cpsUA)(cpsUB)(cpsUC)(cpsUD)(cpsUE)(cpsUF)(cpsUG)(cpsUH)(cpsUI)(cpsUJ)(cpsUK)(cpsUL)(cpsUM)(cpsUN)(cpsUO)(cpsUP)(cpsUQ)(cpsUR)(cpsUS)(cpsUT)(cpsUU)(cpsUV)(cpsUW)(cpsUX)(cpsUY)(cpsUZ)(cpsVA)(cpsVB)(cpsVC)(cpsVD)(cpsVE)(cpsVF)(cpsVG)(cpsVH)(cpsVI)(cpsVJ)(cpsVK)(cpsVL)(cpsVM)(cpsVN)(cpsVO)(cpsVP)(cpsVQ)(cpsVR)(cpsVS)(cpsVT)(cpsVU)(cpsVV)(cpsVW)(cpsVX)(cpsVY)(cpsVZ)(cpsWA)(cpsWB)(cpsWC)(cpsWD)(cpsWE)(cpsWF)(cpsWG)(cpsWH)(cpsWI)(cpsWJ)(cpsWK)(cpsWL)(cpsWM)(cpsWN)(cpsWO)(cpsWP)(cpsWQ)(cpsWR)(cpsWS)(cpsWT)(cpsWU)(cpsWV)(cpsWW)(cpsWX)(cpsWY)(cpsWZ)(cpsXA)(cpsXB)(cpsXC)(cpsXD)(cpsXE)(cpsXF)(cpsXG)(cpsXH)(cpsXI)(cpsXJ)(cpsXK)(cpsXL)(cpsXM)(cpsXN)(cpsXO)(cpsXP)(cpsXQ)(cpsXR)(cpsXS)(cpsXT)(cpsXU)(cpsXV)(cpsXW)(cpsXX)(cpsXY)(cpsXZ)(cpsYA)(cpsYB)(cpsYC)(cpsYD)(cpsYE)(cpsYF)(cpsYG)(cpsYH)(cpsYI)(cpsYJ)(cpsYK)(cpsYL)(cpsYM)(cpsYN)(cpsYO)(cpsYP)(cpsYQ)(cpsYR)(cpsYS)(cpsYT)(cpsYU)(cpsYV)(cpsYW)(cpsYX)(cpsYY)(cpsYZ)(cpsZA)(cpsZB)(cpsZC)(cpsZD)(cpsZE)(cpsZF)(cpsZG)(cpsZH)(cpsZI)(cpsZJ)(cpsZK)(cpsZL)(cpsZM)(cpsZN)(cpsZO)(cpsZP)(cpsZQ)(cpsZR)(cpsZS)(cpsZT)(cpsZU)(cpsZV)(cpsZW)(cpsZX)(cpsZY)(cpsZZ)	83	750	750





TABLE I

S. pneumoniae - Coding regions containing known sequences

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match gene name	Percent Ident	HSP nt length	ORF nt length
41	13	9233	9132	emb127775 SP15	S.pneumoniae DNA for insertion sequence IS1381 (1966 bp)	95	160	402
41	14	9639	9475	emb128200 SP28	S.pneumoniae pcqA gene and open reading frames	100	189	195
44	5	7190	7555	emb128200 SP28	S.pneumoniae pcqA gene and open reading frames	99	366	366
44	6	8059	7607	emb127775 SP15	S.pneumoniae DNA for insertion sequence IS1318 (1172 bp)	97	453	453
44	7	8023	8022	emb127775 SP15	S.pneumoniae DNA for insertion sequence IS1381 (1966 bp)	95	160	402
44	8	8559	8365	emb128200 SP28	S.pneumoniae pcqA gene and open reading frames	100	189	195
48	9	6480	4687	gb U39074	Streptococcus pneumoniae pyruvate oxidase (pxoA) gene, complete cds	99	1794	1794
49	2	231	2603	gb U20561	Streptococcus pneumoniae xpf gene, partial cds	100	216	2373
53	6	2407	2156	gb U04047	Streptococcus pneumoniae 552 dextran glucosidase gene and insertion sequence IS1202 transposase gene, complete cds	97	242	252
53	7	2566	2405	emb128333 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and a11A gene	100	94	162
53	8	2831	2475	emb128333 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and a11A gene	99	338	357
54	13	12409	11105	emb128333 SP28	S.pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamose biosynthesis genes and a11A gene	67	591	1305
55	22	20488	19949	emb1284379 SP28	S.pneumoniae dfr gene (isolate 97)	99	540	540
61	11	11864	5900	emb1216082 NAL	Streptococcus pneumoniae a11B gene	98	1945	1965
63	1	3	239	gb U18729	S.pneumoniae a11A repair protein (hexA) gene, complete cds	100	237	237
63	2	233	2611	gb U18729	S.pneumoniae a11A repair protein (hexA) gene, complete cds	99	2330	2379
63	3	2557	2823	gb U18729	S.pneumoniae a11A repair protein (hexA) gene, complete cds	99	246	267
63	4	2958	4664	gb U18729	S.pneumoniae a11A repair protein (hexA) gene, complete cds	95	69	1707
67	6	3770	3399	gb U20570	Streptococcus pneumoniae hyaluronidase gene, complete cds	96	372	372
67	7	7161	4171	gb U20570	Streptococcus pneumoniae hyaluronidase gene, complete cds	99	2938	2991
70	1	1	702	gb U14140	S.pneumoniae dpm gene region encoding dpmC and dpmD, complete cds	100	693	702
70	2	678	1160	gb U14140	S.pneumoniae dpm1 gene region encoding dpmC and dpmD, complete cds	100	483	483
70	3	2490	1210	gb U14339	S.pneumoniae ftnII gene region encoding ftnII and ftnIII, complete cds	98	462	1281
70	7	4230	4424	gb U04234	S.pneumoniae exodeoxyribonuclease (exoA) gene, complete cds	99	147	395
70	8	5197	4316	gb U04234	S.pneumoniae exodeoxyribonuclease (exoA) gene, complete cds	99	881	882

S. pneumoniae - Coding regions containing known sequences

TABLE 1

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	Percent ident	ORF length	ORF nt
70	13	8108	9874	gb U20562	S. pneumoniae tppd gene, partial cds	93	234	1767
71	22	27964	28341	emb X63602 EP80	S. pneumoniae masA-box	93	233	177
72	5	4607	3552	emb 226850 EPAT	S. pneumoniae (M22) genes for ATPase a subunit, ATPase b subunit and ATPase c subunit	97	102	1056
73	1	471	133	emb X63602 EP80	S. pneumoniae masA-box	91	193	339
73	3	3658	977	gb U30479	S. pneumoniae DNA polymerase I (polA) gene, complete cds	99	2682	2682
73	8	4864	5379	gb H36180	Streptococcus pneumoniae transposase, (comA and comB) and SAICAR synthetase (purC) genes, complete cds	98	318	516
77	3	2622	1999	emb 281335 EP28	S. pneumoniae dexB, cap1A, B, C, D, E, F, G, H, I, J, K, L biosynthesis genes and aIIa gene	95	624	624
77	4	3341	2523	emb 281335 EP28	S. pneumoniae dexB, cap1A, B, C, D, E, F, G, H, I, J, K, L biosynthesis genes and aIIa gene	91	819	819
78	1	341	3	emb X77249 EP6	S. pneumoniae (R6) clbA/clbH genes	99	339	339
78	2	1095	325	emb X77249 EP6	S. pneumoniae (R6) clbA/clbH genes	99	771	771
82	10	11438	10616	gb U90721	Streptococcus pneumoniae signal peptidase I (spi) gene, complete cds	97	621	621
82	11	12402	11134	gb U93576	Streptococcus pneumoniae ribonuclease H1 (rnhA) gene, complete cds	98	953	969
82	12	12381	12704	gb U93576	Streptococcus pneumoniae ribonuclease H1 (rnhB) gene, complete cds	100	51	324
83	8	3212	3550	emb 277727 EP25	S. pneumoniae DNA for insertion sequence IS1318 (822 bp)	97	280	339
83	10	4662	6651	gb H36180	Streptococcus pneumoniae transposase, (comA and comB) and SAICAR synthetase (purC) genes, complete cds	99	2190	2190
83	11	6849	8213	gb H36180	Streptococcus pneumoniae transposase, (comA and comB) and SAICAR synthetase (purC) genes, complete cds	99	1305	1305
83	12	8236	9090	gb H36180	Streptococcus pneumoniae transposase, (comA and comB) and SAICAR synthetase (purC) genes, complete cds	99	955	855
83	13	9283	13017	gb L15190	Streptococcus pneumoniae SAICAR synthetase (purC) gene, complete cds	100	107	3735
83	23	22147	23313	gb L36923	Streptococcus pneumoniae beta-N-acetylhexosaminidase (scrH) gene, complete cds	98	218	1167
83	24	23268	23350	gb L36923	Streptococcus pneumoniae beta-N-acetylhexosaminidase (scrH) gene, complete cds	98	172	183
83	25	27527	23505	gb L36923	Streptococcus pneumoniae beta-N-acetylhexosaminidase (scrH) gene, complete cds	99	1826	4023

TABLE 1

S. pneumoniae - Coding regions containing known sequences

Contig ID	Gen ID	Start (nt)	Stop (nt)	Match	Match gene name	Percent Ident	HSP nt length	ORF nt length
83	26	28472	27771	gb L35923	Streptococcus pneumoniae beta-n-acetylthioaminoimidase (atm) gene, complete cds	99	416	702
84	4	4554	6173	emb 283335 sp28	S.pneumoniae dexB, cap1A,B,C,D,E,F,G,H,I,J,K genes, dTDP-ribose biosynthesis genes and a11a gene	98	697	1620
85	6	5531	5216	emb 277725 sp28	S.pneumoniae RNA for insertion sequence IS1381 (966 bp)	96	439	636
88	5	2357	3511	gb M36180	Streptococcus pneumoniae transposase, (conA and comB) and SAIAR synthetase (parC) genes, complete cds	94	555	555
88	6	3466	4269	gb M36180	Streptococcus pneumoniae transposase, (conA and comB) and SAIAR synthetase (parC) genes, complete cds	94	804	804
89	13	9778	10093	gb M36180	Streptococcus pneumoniae transposase, (conA and comB) and SAIAR synthetase (parC) genes, complete cds	97	211	216
89	14	10662	10412	emb 283335 sp28	S.pneumoniae dexB, cap1A,B,C,D,E,F,G,H,I,J,K genes, dTDP-ribose biosynthesis genes and a11a gene	97	335	351
93	10	5303	4941	emb M63602 sp80	S.pneumoniae swaB-box	89	237	363
97	4	1708	1520	gb M41735	Streptococcus pneumoniae peptide methionine sulfoxide reductase (msrA) and homoserine kinase homolog (thrB) genes, complete cds	91	140	189
99	1	89	700	emb 283335 sp28	S.pneumoniae dexB, cap1A,B,C,D,E,F,G,H,I,J,K genes, dTDP-ribose biosynthesis genes and a11a gene	93	592	612
99	2	1773	775	emb M17337 spAM	Streptococcus pneumoniae aml locus conferring aminopterin resistance	99	998	999
99	3	2794	1712	emb M17337 spAM	Streptococcus pneumoniae aml locus conferring aminopterin resistance	99	1083	1083
99	4	3732	2788	emb M17337 spAM	Streptococcus pneumoniae aml locus conferring aminopterin resistance	100	945	945
99	5	5249	3714	emb M17337 spAM	Streptococcus pneumoniae aml locus conferring aminopterin resistance	100	1536	1536
99	6	7262	5277	emb M17337 spAM	Streptococcus pneumoniae aml locus conferring aminopterin resistance	99	1986	1986
101	1	216	1538	emb M54225 spEN	S.pneumoniae epuA and euh genes for 7 kDa protein and membrane endonuclease	99	146	1323
101	2	1492	1719	emb M54225 spEN	S.pneumoniae epuA k12 euh genes for 7 kDa protein and membrane endonuclease	99	228	228
101	3	1694	1855	emb M54225 spEN	S.pneumoniae epuA and euh genes for 7 kDa protein and membrane endonuclease	100	162	162
101	4	1701	2882	emb M54225 spEN	S.pneumoniae epuA and euh genes for 7 kDa protein and membrane endonuclease	100	882	882
103	7	2556	5041	emb 295914 sp29	Streptococcus pneumoniae sodA gene	100	396	516
104	2	1347	1556	emb 277727 spIS	S.pneumoniae RNA for insertion sequence IS1318 (823 bp)	83	206	210

S. pneumoniae - Coding regions containing known sequences

TABLE 1

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	Accession ident	ORF nt length	ORF nt length
105	5	5381	5028	emb 567733 SPPA	S. pneumoniae parC, parE and transposase genes and unknown orf	98	353	354
105	6	5049	5319	emb 567733 SPPA	S. pneumoniae parC, parE and transposase genes and unknown orf	98	84	711
107	4	2785	1880	emb X11022 SPPE	S. pneumoniae penA gene	98	72	906
107	5	2913	4388	emb X11022 SPPE	S. pneumoniae penA gene	99	1692	2076
107	6	4981	5395	emb X11136 SPPE	Streptococcus pneumoniae penA gene for penicillin binding protein 2B lacking N-term. (penicillin resistant strain)	91	107	615
108	9	9058	8718	emb 567733 SPPA	S. pneumoniae parC, parE and transposase genes and unknown orf	95	342	351
108	12	11308	10522	emb 567733 SPPA	S. pneumoniae parC, parE and transposase genes and unknown orf	99	199	387
109	3	2768	2241	emb 577725 SPIS	S. pneumoniae DNA for insertion sequence IS1381 (965 bp)	96	61	528
109	4	2688	2855	emb 577725 SPIS	S. pneumoniae DNA for insertion sequence IS1318 (1372 bp)	96	148	168
109	5	2862	3269	emb 577727 SPIS	S. pneumoniae DNA for insertion sequence IS1318 (823 bp)	97	353	408
109	6	5320	3584	gb M58729	S. pneumoniae mismatch repair protein (hesA) gene, complete cds	100	371	1737
113	1	431	3	gb M56180	Streptococcus pneumoniae transposase, (cmaA and cmaB) and SATCAR synthetase (purC) genes, complete cds	95	429	429
113	2	9788	8532	emb 593400 SPDA	S. pneumoniae dnaA gene and orf	99	1257	1257
113	11	9870	10985	emb 593400 SPDA	S. pneumoniae dnaA gene and orf	99	1116	1116
114	3	2530	2030	gb M56180	Streptococcus pneumoniae transposase, (cmaA and cmaB) and SATCAR synthetase (purC) genes, complete cds	95	481	501
115	11	11303	10932	gb U01047	Streptococcus pneumoniae SZ2 dextran glucosidase gene and insertion sequence IS1202 transposase gene, complete cds	97	372	372
117	1	897	3102	emb X72967 SPNA	S. pneumoniae nanA gene	99	2402	2406
117	2	3277	3331	emb X72967 SPNA	S. pneumoniae nanA gene	99	237	555
117	3	4327	3899	gb M56180	Streptococcus pneumoniae transposase, (cmaA and cmaB) and SATCAR synthetase (purC) genes, complete cds	98	429	429
121	2	1369	1341	gb U01210	Streptococcus pneumoniae heat shock protein 70 (dnaK) gene, complete cds and dnaJ (dnaJ) gene, partial cds	99	202	573
121	3	2412	4353	gb U01210	Streptococcus pneumoniae heat shock protein 70 (dnaK) gene, complete cds and dnaJ (dnaJ) gene, partial cds	99	1842	1842
122	8	5066	5587	gb U01047	Streptococcus pneumoniae SZ2 dextran glucosidase gene and insertion sequence IS1202 transposase gene, complete cds	64	451	522

S. pneumoniae - Coding regions containing known sequences

TABLE 1

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	percent ident	RSP nt length	ORF nt length
125	1	1811	189	[gb]M3180	Streptococcus pneumoniae transposase, (com and comB) and SATCM synthetase (purC) genes, complete cds	92	99	1623
126	15	12496	11204	[emb]J23335 SP28	S. pneumoniae dexB, cap1A,B,C,D,E,F,G,H,I,J,K genes, dTDP-thiamose operon, and a1a gene	91	705	1253
134	1	1	492	[emb]Y10818 SPV1	S. pneumoniae spa gene	99	203	492
134	2	556	2452	[gb]M319304	Streptococcus pneumoniae choline binding protein A (cbpA) gene, partial cds	86	685	2097
134	3	1160	837	[emb]Y10818 SPV1	S. pneumoniae spa gene	86	324	324
134	4	3392	2882	[gb]M319304	Streptococcus pneumoniae choline binding protein A (cbpA) gene, partial cds	98	215	1071
134	8	7992	9848	[gb]J12567	Streptococcus pneumoniae pil3 glycerol-3-phosphate dehydrogenase (glpD) gene, partial cds, and glycerol uptake facilitator (glpF) and ORF3 genes, complete cds	99	285	1857
134	9	9846	10622	[gb]J12567	Streptococcus pneumoniae pil3 glycerol-3-phosphate dehydrogenase (glpD) gene, partial cds, and glycerol uptake facilitator (glpF) and ORF3 genes, complete cds	99	570	777
134	10	10605	11122	[gb]J12567	Streptococcus pneumoniae pil3 glycerol-3-phosphate dehydrogenase (glpD) gene, partial cds, and glycerol uptake facilitator (glpF) and ORF3 genes, complete cds	100	318	318
137	13	7970	8443	[gb]J09239	Streptococcus pneumoniae type 139 capsular polysaccharide biosynthesis operon, (cpsA BCDEF H I J K L) genes, complete cds, and a1a gene, partial cds	90	420	474
137	14	8250	8775	[emb]283335 SP28	S. pneumoniae dexB, cap1A,B,C,D,E,F,G,H,I,J,K genes, dTDP-thiamose biosynthesis genes and a1a gene	94	174	186
137	15	8773	8967	[emb]283335 SP28	S. pneumoniae dexB, cap1A,B,C,D,E,F,G,H,I,J,K genes, dTDP-thiamose biosynthesis genes and a1a gene	98	195	195
137	16	9222	9667	[emb]27726 SP15	S. pneumoniae DNA for insertion sequence IS1318 (1372 bp)	96	446	465
137	17	9641	10051	[emb]27727 SP15	S. pneumoniae DNA for insertion sequence IS1318 (823 bp)	96	293	411
139	10	12398	12702	[emb]M33602 SP80	S. pneumoniae msaA-Box	90	234	297
141	8	7805	8938	[emb]249988 SPM	Streptococcus pneumoniae msaA gene	99	338	1134
141	9	8936	10972	[emb]249988 SPM	Streptococcus pneumoniae msaA gene	99	2037	2017
141	10	11472	12467	[emb]249988 SPM	Streptococcus pneumoniae msaA gene	100	76	996
142	2	257	814	[gb]M80215	Streptococcus pneumoniae uvs402 protein gene, complete cds	98	174	558
142	3	787	957	[gb]M80215	Streptococcus pneumoniae uvs402 protein gene, complete cds	100	142	171
142	4	980	3022	[gb]M80215	Streptococcus pneumoniae uvs402 protein gene, complete cds	95	197	2603

TABLE 1

S. pneumoniae - Coding regions containing known sequences

Contig ID	Start (nt)	Stop (nt)	Accession	Match gene name	percent ident	HSP nt length	ORF nt length
142	5	3020	gb H80215	Streptococcus pneumoniae uva402 protein gene, complete cds	100	153	576
145	1	219	emb Z23513 SPAL	S. pneumoniae p13a gene for anti-late gene A	97	185	219
145	2	171	gb L20556	Streptococcus pneumoniae plpA gene, partial cds	99	181	184
145	3	2287	gb Z47210 SPR6	S. pneumoniae dead, capB, capB and capC genes and orfA	99	1052	5313
145	4	9314	gb H90527	Streptococcus pneumoniae penicillin-binding protein (penA) gene, complete cds	99	2165	2169
145	5	10488	gb H90527	Streptococcus pneumoniae penicillin-binding protein (penA) gene, complete cds	99	512	567
146	1	139	emb Z23002 SP28	S. pneumoniae pcgA and pcgC genes	98	156	156
146	2	344	emb Z23002 SP28	S. pneumoniae pcgB and pcgC genes	98	255	255
146	16	11795	gb H31801	S. pneumoniae pcgB and pcgC genes	85	276	1052
147	11	10578	emb Z21702 SPUN	S. pneumoniae ung gene and mutX genes encoding uracil-DNA glycosylase and 8-oxodGTP nucleoside triphosphatase	98	477	477
147	12	11338	emb Z21702 SPUN	S. pneumoniae ung gene and mutX genes encoding uracil-DNA glycosylase and 8-oxodGTP nucleoside triphosphatase	99	663	663
148	12	9009	gb U47131	Streptococcus pneumoniae peptide methionine sulfoxide reductase (metR) and methionine sulfoxide lyase (metS) genes, complete cds	90	180	195
156	4	1154	emb X63602 SPBO	S. pneumoniae msaA-Box	94	185	249
159	13	9048	gb H31801	Streptococcus pneumoniae msaA-Box	98	526	528
160	1	1	emb Z26851 SPAT	S. pneumoniae (R6) genes for ATPase a subunit, ATPase b subunit and ATPase c subunit	100	142	147
160	2	179	emb Z26851 SPAT	S. pneumoniae (R6) genes for ATPase a subunit, ATPase b subunit and ATPase c subunit	99	720	720
160	3	906	emb Z26850 SPAT	S. pneumoniae (H222) genes for ATPase a subunit, ATPase b subunit and ATPase c subunit	95	501	501
160	4	1373	emb Z26850 SPAT	S. pneumoniae (H222) genes for ATPase a subunit, ATPase b subunit and ATPase c subunit	87	306	570
161	1	1	emb X77249 SPK6	S. pneumoniae (R6) ciar/ciaH genes	99	984	984
161	7	6910	emb X83917 SPGY	S. pneumoniae orfgyrA and gyrB gene encoding DNA gyrase B subunit	99	417	568
161	8	7443	emb X83917 SPGY	S. pneumoniae orfgyrA and gyrB gene encoding DNA gyrase B subunit	96	192	1944
161	1	2	emb X83917 SPGY	S. pneumoniae orfgyrA and gyrB gene encoding DNA gyrase B subunit	98	327	2154

TABLE I

S. pneumoniae - Coding regions containing known sequences

GenBank ID	Start ID	Stop ID	match accession	match gene name	percent ident	HSP nt length	ORF nt length
165	1	32	gb J01786	S.pneumoniae malX and malM genes encoding membrane protein and arylsulphatase, complete cds, and malP gene encoding phosphorylase	93	1587	1587
165	2	1608	gb J01786	S.pneumoniae malX and malM genes encoding membrane protein and arylsulphatase, complete cds, and malP gene encoding phosphorylase	100	280	2295
166	1	378	emb J11463 SP01	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	100	375	375
166	2	1507	320	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	99	1188	1188
166	3	3240	1412	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	99	563	1809
167	1	1077	328	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	94	155	750
167	2	1844	959	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	98	405	846
167	3	2714	1842	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	97	604	873
167	4	3399	2641	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	99	703	759
168	1	1	2359	Streptococcus pneumoniae dnaG, rpoD, cpkA genes and ORF3 and ORF5	99	282	2259
170	10	7138	7645	S.pneumoniae DNA for insertion sequence IS118 (1172 bp)	95	315	318
172	6	2462	4981	Streptococcus pneumoniae formate acetyltransferase (exp72) gene, partial cds	97	365	2520
175	1	373	20	Streptococcus pneumoniae transposase, [comA and comB] and SAICOM synthetase (exp72) gene, complete cds	89	353	354
175	4	1843	3621	S.pneumoniae dexA, cap3A, cap3B and cap3C genes and orfs	95	89	1779
176	5	3984	2980	S.pneumoniae parC, parE and transposase genes and unknown orf	100	573	1005
178	1	3	425	S.pneumoniae parC, parE and transposase genes and unknown orf	95	423	423
179	1	426	70	S.pneumoniae dexA, cap3A, B.C.D.E.F.G.H.I.J.K.L genes, dTDP-thiamose synthetase gene	99	338	357
180	3	3084	1855	S.pneumoniae gyrA gene	99	381	1220
186	1	714	4	emb J27691 SD08 S.pneumoniae ynfJ, B.C.D.E.I, ftsL, plxP and regH genes	98	59	711
186	2	2254	408	S.pneumoniae ynfJ, B.C.D.E.I, ftsL, plxP and regH genes	98	315	1867
186	3	707	800	S.pneumoniae ynfJ, B.C.D.E.I, ftsL, plxP and regH genes	98	174	174
189	1	2	259	Streptococcus pneumoniae heat shock protein 70 (dnaK) gene, complete cds and dnaJ (dnaJ) gene, partial cds	99	258	258
189	2	600	385	Streptococcus pneumoniae heat shock protein 70 (dnaK) gene, complete cds and dnaJ (dnaJ) gene, partial cds	98	204	216

TABLE 1

S. pneumoniae - Coding regions containing known sequences

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	percent match	HSP RE	ORF RE	ORF length
189	3	1018	851	gb U72720	Streptococcus pneumoniae heat shock protein 70 (dnaK) gene, complete cds and dnaK (dnaK) gene, partial cds	99	158	168	
189	4	1012	2134	gb U72720	Streptococcus pneumoniae heat shock protein 70 (dnaK) gene, complete cds and dnaK (dnaK) gene, partial cds	99	1062	1143	
191	9	7829	7524	emb X63602 SPEO	S. pneumoniae mask-Box				
194	1	1	725	gb U61810	Streptococcus pneumoniae transposase, (comA and comB) and SAICAR synthetase (purC) genes, complete cds	95	214	306	
199	2	1117	881	emb Z83335 SP28	S. pneumoniae dexA, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	91	728	729	
199	4	1499	1762	emb Z83335 SP28	S. pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	96	211	237	
199	5	1781	2284	emb Z83335 SP28	S. pneumoniae dexC, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	89	248	264	
203	1	1277	1277	gb Z20543	Streptococcus pneumoniae Exp gene, partial cds	98	504	504	
204	1	1145	3	gb L16131	Streptococcus pneumoniae exp10 gene, complete cds	99	342	1643	
208	1	59	2296	gb U69711	Streptococcus pneumoniae pneumococcal surface protein A PspA (pspA) gene, partial cds	99	1143	1143	
213	3	2455	2123	emb Z83335 SP28	S. pneumoniae dexA, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	96	332	333	
216	1	368	12	emb Z83335 SP28	S. pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	99	338	357	
216	3	2650	2377	gb U48678	S. pneumoniae promoter sequence DNA	98	86	324	
222	1	417	4	emb Z83335 SP28	S. pneumoniae dexA, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	94	414	414	
227	3	5266	4238	emb A2000336 SP	Streptococcus pneumoniae lch gene	99	1029	1039	
239	1	1	804	gb H31266	S. pneumoniae recP gene, complete cds	95	484	804	
247	3	1625	1807	gb H51510	Streptococcus pneumoniae transposase, (comA and comB) and SAICAR synthetase (purC) genes, complete cds	94	178	183	
249	3	921	1164	emb Z83335 SP28	S. pneumoniae dexB, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	94	413	444	
253	1	362	3	gb H31810	Streptococcus pneumoniae transposase, (comA and comB) and SAICAR synthetase (purC) genes, complete cds	99	360	360	
253	5	1238	2050	emb Z83335 SP28	S. pneumoniae dexA, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, RTDP-rhamnose biosynthesis genes and allA gene	95	420	813	



TABLE 1  
S. pneumoniae - Coding regions containing known sequences

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match gene name	Percent ident	MSP nt length	ORF nt length
253	6	2069	2572	emb 263115 SP26	S.pneumoniae dssA, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamase biosynthesis genes and allA gene	97	504	504
255	1	3	800	emb 262002 SP26	S.pneumoniae popA and popC genes	97	531	798
255	2	798	1841	emb 262002 SP26	S.pneumoniae popA and popC genes	97	672	1044
255	3	2493	1369	emb 267739 SPPA	S.pneumoniae parC, parE and ... genes and unknown orf	92	435	535
257	2	985	770	emb K17337 SPM	Streptococcus pneumoniae ami locus conferring aminopterin resistance	96	117	216
257	3	1245	907	gb H56180	Streptococcus pneumoniae transposase, (comb and comb) and SAICAR synthetase [purC] genes, complete cds	97	339	339
267	2	495	1208	gb U14136	Streptococcus pneumoniae dihydropterate synthase (sulM), dihydrofolate synthetase (sulB), guanosine triphosphate cyclohydrolase (sulC), aldolase-pyrophosphokinase (sulD) genes, complete cds	95	84	714
267	3	1291	2277	gb U14136	Streptococcus pneumoniae dihydropterate synthase (sulM), dihydrofolate synthetase (sulB), guanosine triphosphate cyclohydrolase (sulC), aldolase-pyrophosphokinase (sulD) genes, complete cds	97	755	927
267	4	2261	3601	gb U14136	Streptococcus pneumoniae dihydropterate synthase (sulM), dihydrofolate synthetase (sulB), guanosine triphosphate cyclohydrolase (sulC), aldolase-pyrophosphokinase (sulD) genes, complete cds	98	1341	1341
267	5	3581	4136	gb U14136	Streptococcus pneumoniae dihydropterate synthase (sulM), dihydrofolate synthetase (sulB), guanosine triphosphate cyclohydrolase (sulC), aldolase-pyrophosphokinase (sulD) genes, complete cds	99	576	576
267	6	4164	4949	gb U14136	Streptococcus pneumoniae dihydropterate synthase (sulM), dihydrofolate synthetase (sulB), guanosine triphosphate cyclohydrolase (sulC), aldolase-pyrophosphokinase (sulD) genes, complete cds	99	748	786
267	7	5344	5140	gb U14136	Streptococcus pneumoniae dihydropterate synthase (sulM), dihydrofolate synthetase (sulB), guanosine triphosphate cyclohydrolase (sulC), aldolase-pyrophosphokinase (sulD) genes, complete cds	100	186	405
286	4	1793	1990	emb K63602 SPRO	S.pneumoniae asacA-Box			
271	1	562	104	gb J029486	S.pneumoniae ... patch repair (hexB) genes, complete cds	89	194	198
291	1	75	524	gb J04047	Streptococcus pneumoniae S22 dactin glucosidase gene and insertion sequence IS1207 transposase gene, complete cds	93	160	459
291	2	1001	525	emb 263335 SP26	S.pneumoniae dssA, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamase biosynthesis genes and allA gene	96	450	450
291	3	807	559	emb 263335 SP26	S.pneumoniae dssA, cap1(A,B,C,D,E,F,G,H,I,J,K) genes, dTDP-thiamase biosynthesis genes and allA gene	87	205	477
291	4	1374	1099	gb H16180	Streptococcus pneumoniae transposase, (comb and comb) and SAICAR synthetase [purC] genes, complete cds	90	170	249
					Streptococcus pneumoniae transposase, (comb and comb) and SAICAR synthetase [purC] genes, complete cds	85	244	276

TABLE 1  
S. pneumoniae Coding regions containing known sequences

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	match gene name	percent ident	HSP nt length	ORF nt length
233	1	3	1673	emb 267740 SP07	S. pneumoniae gyrB gene and unknown orf	98	533	1671
286	1	1434	151	emb 247210 SP16	S. pneumoniae dexB, capsA, capsB and capsC genes and orfs	99	430	1284
317	1	157	510	emb 267739 SP14	S. pneumoniae parC, parE and transposase genes and unknown orf	89	353	354
325	2	1237	485	emb 283335 SP28	S. pneumoniae capsA, C, D, E, F, G, H, I, J, K genes, dtdp-rhamnose biosynthesis genes and orfs	91	299	753
326	1	667	102	emb 228201 SP28	S. pneumoniae pcpsA gene and open reading frames	100	233	462
327	1	603	64	emb 283335 SP28	S. pneumoniae dexA, capsA, B, C, D, E, F, G, H, I, J, K genes, QTP-rhamnose biosynthesis genes and a11a gene	94	89	540
334	1	153	545	gb 041735	S. pneumoniae pneumolysin, pneumolysin peptide methionine sulfoxide reductase (narA) and homoserine V	87	91	393
336	1	308	93	emb 226650 SPMT	S. pneumoniae (H222) genes for ATPase & subunit, ATPase b subunit and ATPase c subunit	97	102	216
360	1	1	519	emb 267739 SP14	S. pneumoniae parC, parE and transposase genes and unknown orf	95	435	519
360	4	1398	1960	emb 283335 SP28	S. pneumoniae dexB, capsA, B, C, D, E, F, G, H, I, J, K genes, dtdp-rhamnose biosynthesis genes and a11a gene	94	353	363
362	1	673	2	emb 283335 SP28	S. pneumoniae dexB, capsA, B, C, D, E, F, G, H, I, J, K genes, dtdp-rhamnose biosynthesis genes and a11a gene	95	63	672
362	2	1168	728	gb 040471	Streptococcus pneumoniae SS2 dectran adherence gene and insertion sequence IS1202 transposase gene, complete cds	96	441	441
384	1	347	111	emb 285787 SPCP	S. pneumoniae capsA, capsB, capsC, capsD, capsE, capsF, capsG, capsH, capsI, capsJ, capsK, capsL, capsM, capsN, capsO, capsP, capsQ, capsR, capsS, capsT, capsU, capsV, capsW, capsX, capsY, capsZ, capsAA, capsAB, capsAC, capsAD, capsAE, capsAF, capsAG, capsAH, capsAI, capsAJ, capsAK, capsAL, capsAM, capsAN, capsAO, capsAP, capsAQ, capsAR, capsAS, capsAT, capsAU, capsAV, capsAW, capsAX, capsAY, capsAZ, capsBA, capsBB, capsBC, capsBD, capsBE, capsBF, capsBG, capsBH, capsBI, capsBJ, capsBK, capsBL, capsBM, capsBN, capsBO, capsBP, capsBQ, capsBR, capsBS, capsBT, capsBU, capsBV, capsBW, capsBX, capsBY, capsBZ, capsCA, capsCB, capsCC, capsCD, capsCE, capsCF, capsCG, capsCH, capsCI, capsCJ, capsCK, capsCL, capsCM, capsCN, capsCO, capsCP, capsCQ, capsCR, capsCS, capsCT, capsCU, capsCV, capsCW, capsCX, capsCY, capsCZ, capsDA, capsDB, capsDC, capsDD, capsDE, capsDF, capsDG, capsDH, capsDI, capsDJ, capsDK, capsDL, capsDM, capsDN, capsDO, capsDP, capsDQ, capsDR, capsDS, capsDT, capsDU, capsDV, capsDW, capsDX, capsDY, capsDZ, capsEA, capsEB, capsEC, capsED, capsEE, capsEF, capsEG, capsEH, capsEI, capsEJ, capsEK, capsEL, capsEM, capsEN, capsEO, capsEP, capsEQ, capsER, capsES, capsET, capsEU, capsEV, capsEW, capsEX, capsEY, capsEZ, capsFA, capsFB, capsFC, capsFD, capsFE, capsFF, capsFG, capsFH, capsFI, capsFJ, capsFK, capsFL, capsFM, capsFN, capsFO, capsFP, capsFQ, capsFR, capsFS, capsFT, capsFU, capsFV, capsFW, capsFX, capsFY, capsFZ, capsGA, capsGB, capsGC, capsGD, capsGE, capsGF, capsGG, capsGH, capsGI, capsGJ, capsGK, capsGL, capsGM, capsGN, capsGO, capsGP, capsGQ, capsGR, capsGS, capsGT, capsGU, capsGV, capsGW, capsGX, capsGY, capsGZ, capsHA, capsHB, capsHC, capsHD, capsHE, capsHF, capsHG, capsHH, capsHI, capsHJ, capsHK, capsHL, capsHM, capsHN, capsHO, capsHP, capsHQ, capsHR, capsHS, capsHT, capsHU, capsHV, capsHW, capsHX, capsHY, capsHZ, capsIA, capsIB, capsIC, capsID, capsIE, capsIF, capsIG, capsIH, capsII, capsIJ, capsIK, capsIL, capsIM, capsIN, capsIO, capsIP, capsIQ, capsIR, capsIS, capsIT, capsIU, capsIV, capsIW, capsIX, capsIY, capsIZ, capsJA, capsJB, capsJC, capsJD, capsJE, capsJF, capsJG, capsJH, capsJI, capsJJ, capsJK, capsJL, capsJM, capsJN, capsJO, capsJP, capsJQ, capsJR, capsJS, capsJT, capsJU, capsJV, capsJW, capsJX, capsJY, capsJZ, capsKA, capsKB, capsKC, capsKD, capsKE, capsKF, capsKG, capsKH, capsKI, capsKJ, capsKL, capsKM, capsKN, capsKO, capsKP, capsKQ, capsKR, capsKS, capsKT, capsKU, capsKV, capsKW, capsKX, capsKY, capsKZ, capsLA, capsLB, capsLC, capsLD, capsLE, capsLF, capsLG, capsLH, capsLI, capsLJ, capsLK, capsLL, capsLM, capsLN, capsLO, capsLP, capsLQ, capsLR, capsLS, capsLT, capsLU, capsLV, capsLW, capsLX, capsLY, capsLZ, capsMA, capsMB, capsMC, capsMD, capsME, capsMF, capsMG, capsMH, capsMI, capsMJ, capsMK, capsML, capsMN, capsMO, capsMP, capsMQ, capsMR, capsMS, capsMT, capsMU, capsMV, capsMW, capsMX, capsMY, capsMZ, capsNA, capsNB, capsNC, capsND, capsNE, capsNF, capsNG, capsNH, capsNI, capsNJ, capsNK, capsNL, capsNM, capsNO, capsNP, capsNQ, capsNR, capsNS, capsNT, capsNU, capsNV, capsNW, capsNX, capsNY, capsNZ, capsOA, capsOB, capsOC, capsOD, capsOE, capsOF, capsOG, capsOH, capsOI, capsOJ, capsOK, capsOL, capsOM, capsON, capsOO, capsOP, capsOQ, capsOR, capsOS, capsOT, capsOU, capsOV, capsOW, capsOX, capsOY, capsOZ, capsPA, capsPB, capsPC, capsPD, capsPE, capsPF, capsPG, capsPH, capsPI, capsPJ, capsPK, capsPL, capsPM, capsPN, capsPO, capsPP, capsPQ, capsPR, capsPS, capsPT, capsPU, capsPV, capsPW, capsPX, capsPY, capsPZ, capsQA, capsQB, capsQC, capsQD, capsQE, capsQF, capsQG, capsQH, capsQI, capsQJ, capsQK, capsQL, capsQM, capsQN, capsQO, capsQP, capsQQ, capsQR, capsQS, capsQT, capsQU, capsQV, capsQW, capsQX, capsQY, capsQZ, capsRA, capsRB, capsRC, capsRD, capsRE, capsRF, capsRG, capsRH, capsRI, capsRJ, capsRK, capsRL, capsRM, capsRN, capsRO, capsRP, capsRQ, capsRR, capsRS, capsRT, capsRU, capsRV, capsRW, capsRX, capsRY, capsRZ, capsSA, capsSB, capsSC, capsSD, capsSE, capsSF, capsSG, capsSH, capsSI, capsSJ, capsSK, capsSL, capsSM, capsSN, capsSO, capsSP, capsSQ, capsSR, capsSS, capsST, capsSU, capsSV, capsSW, capsSX, capsSY, capsSZ, capsTA, capsTB, capsTC, capsTD, capsTE, capsTF, capsTG, capsTH, capsTI, capsTJ, capsTK, capsTL, capsTM, capsTN, capsTO, capsTP, capsTQ, capsTR, capsTS, capsTT, capsTU, capsTV, capsTW, capsTX, capsTY, capsTZ, capsUA, capsUB, capsUC, capsUD, capsUE, capsUF, capsUG, capsUH, capsUI, capsUJ, capsUK, capsUL, capsUM, capsUN, capsUO, capsUP, capsUQ, capsUR, capsUS, capsUT, capsUU, capsUV, capsUW, capsUX, capsUY, capsUZ, capsVA, capsVB, capsVC, capsVD, capsVE, capsVF, capsVG, capsVH, capsVI, capsVJ, capsVK, capsVL, capsVM, capsVN, capsVO, capsVP, capsVQ, capsVR, capsVS, capsVT, capsVU, capsVV, capsVW, capsVX, capsVY, capsVZ, capsWA, capsWB, capsWC, capsWD, capsWE, capsWF, capsWG, capsWH, capsWI, capsWJ, capsWK, capsWL, capsWM, capsWN, capsWO, capsWP, capsWQ, capsWR, capsWS, capsWT, capsWU, capsWV, capsWW, capsWX, capsWY, capsWZ, capsXA, capsXB, capsXC, capsXD, capsXE, capsXF, capsXG, capsXH, capsXI, capsXJ, capsXK, capsXL, capsXM, capsXN, capsXO, capsXP, capsXQ, capsXR, capsXS, capsXT, capsXU, capsXV, capsXW, capsXX, capsXY, capsXZ, capsYA, capsYB, capsYC, capsYD, capsYE, capsYF, capsYG, capsYH, capsYI, capsYJ, capsYK, capsYL, capsYM, capsYN, capsYO, capsYP, capsYQ, capsYR, capsYS, capsYT, capsYU, capsYV, capsYW, capsYX, capsYY, capsYZ, capsZA, capsZB, capsZC, capsZD, capsZE, capsZF, capsZG, capsZH, capsZI, capsZJ, capsZK, capsZL, capsZM, capsZN, capsZO, capsZP, capsZQ, capsZR, capsZS, capsZT, capsZU, capsZV, capsZW, capsZX, capsZY, capsZZ	94	54	237

5. pneumoniae - Putative coding regions of novel proteins similar to known proteins

TABLE 2

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
228	2	1760	1942	gi P60663 P606	translational elongation factor Tu - Streptococcus oralis	100	100	183
319	1	2	205	gi P94927	neomycin phosphotransferase (cloning vector pBS29)	100	100	204
260	1	2	1138	gi P60663 P606	translational elongation factor Tu - Streptococcus oralis	99	98	1137
25	2	486	1384	gi 1574495	hypothetical [Haemophilus influenzae]	98	96	909
94	2	685	1002	gi 1310627	phosphoenolpyruvate:sugar phosphotransferase system IIF [Streptococcus mutans]	98	93	318
312	1	190	2	gi 1347999	ATP-dependent protease proteolytic subunit [Streptococcus salivarius]	98	95	189
329	1	1	807	gi 924848	inosine monophosphate dehydrogenase [Streptococcus pneumoniae]	96	94	807
316	2	290	589	gi 987050	lacZ gene product [unidentified cloning vector]	98	98	300
181	9	5948	7366	gi 155755	phospho-beta-D-galactosidase (EC 3.2.1.45) [Streptococcus lactis cremoris]	97	94	1419
312	2	1044	361	gi 1347998	uracil phosphoribosyltransferase [Streptococcus salivarius]	97	88	684
32	8	6575	7486	sp P37214 PBA_S	GTP-BINDING PROTEIN BBA KPWLOO	96	91	912
94	3	951	2741	gi 153615	phosphoenolpyruvate:sugar phosphotransferase system enzyme I [Streptococcus salivarius]	96	92	1791
127	1	1	168	gi 581299	initiation factor IF-1 [Lactococcus lactis]	96	90	168
128	14	10438	11154	gi 1278873	deoD [Streptococcus thermophilus]	96	93	717
181	4	1362	1598	gi 46606	lacD polypeptide (AA 1-326) [Staphylococcus aureus]	96	80	237
218	1	1	834	gi 1743856	intragenomic coaggregation-relevant adhesin [Streptococcus gordonii]	96	93	834
319	2	115	441	gi 208225	heat-shock protein 82/neomycin phosphotransferase (cloning vector)	96	96	327
54	12	8622	10867	gi P0100972	pyruvate formate-lyase [Streptococcus mutans]	95	99	2346
181	2	666	1289	gi 149396	lacD [Lactococcus lactis]	95	80	684
45	3	3405	3605	gi 1350606	fixn [Streptococcus mutans]	94	86	366
89	10	7972	7337	gi 703442	thymidine kinase [Streptococcus gordonii]	94	80	636
148	9	6201	7354	gi 159767	[ORF-glucose pyrophosphorylase [Streptococcus pyogenes]]	94	85	924
160	7	4430	5848	gi 151573	[ORF-ATPase [Enterococcus faecalis]]	94	97	1419
2	2	4598	3513	gi 151763	plasma receptor [Streptococcus pyogenes]	93	86	1086
12	8	7877	6204	gi 1103845	formyl-tetrahydrofolate synthetase [Streptococcus mutans]	93	84	1674

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

TABLE 2

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
65	11	4734	5120	gi 40150	Lt4 protein (AX 1-122) [Bacillus subtilis]	93	87	387
68	1	53	1297	gi 47341	antitumor protein [Streptococcus pyogenes]	93	87	1245
80	1	3	299	gn JPDJ01166	ribosomal protein S7 [Bacillus subtilis]	93	84	297
127	3	695	1093	gi 104262	ribosomal protein S11 [Bacillus subtilis]	93	86	399
140	5	1924	3462	gi 177326	ATP-binding subunit [Streptococcus mutans]	93	85	1539
211	5	1757	3047	gi 535273	aminopeptidase C [Streptococcus thermophilus]	93	82	711
262	1	16	564	gi 141934	ribosomal protein L2 [Lactococcus lactis]	93	90	549
366	1	197	3	gi 295259	tryptophan synthase beta subunit [Synecococcus sp.]	93	91	195
25	3	1392	1976	gi 157498	hypothetical [Haemophilus influenzae]	92	80	585
36	21	120781	19937	gi 310642	hydrophobic membrane protein [Streptococcus gordonii]	92	86	855
181	3	1265	1536	gi 402296	Lact [Lactococcus lactis]	92	83	270
181	7	3662	4060	gi 149410	enzyme III [Lactococcus lactis]	92	83	399
32	4	563	2937	gn JPDJ029409	flavonectin-binding protein-like protein A [Streptococcus gordonii]	91	85	1695
46	2	3054	1462	gi 1850607	signal recognition particle Fth [Streptococcus mutans]	91	84	1503
65	10	4442	4726	gi 577485 5178	ribosomal protein S17 - Bacillus sterothomophilus	91	80	285
77	2	260	1900	gi 287871	groEL gene product [Lactococcus lactis]	91	82	1811
84	1	2	2056	gi 871784	Clp-like ATP-dependent protease binding subunit [Bos taurus]	91	79	2055
99	8	10750	9272	gi 153740	sucrose phosphorylase [Streptococcus mutans]	91	84	1479
99	9	11947	11072	gi 153739	membrane protein [Streptococcus mutans]	91	78	876
127	5	2065	2469	gi 507223 8585	ribosomal protein L17 - Bacillus sterothomophilus	91	78	405
132	6	9559	9390	gi 143065	hustat [Bacillus sterothomophilus]	91	89	150
137	8	4745	6153	gn JPDJ010047	hex - ATPase beta subunit [Haemophilus influenzae]	91	79	1389
151	7	11119	9734	gi 1815534	glutamine synthetase type 1 [Streptococcus agalactiae]	91	82	1386
201	2	1798	278	gi 2208998	dextran glucosylase [Bacillus sterothomophilus]	91	79	1521
222	2	673	1839	gi 153741	ATP-binding protein [Streptococcus mutans]	91	85	1167
293	5	4113	4400	gi 1156921	unknown protein with function sequence 158611	91	71	288
32	7	6166	6570	gi 1456933 A169	diacylglycerol kinase homolog - Streptococcus mutans	90	77	405

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	Length (nt)
33	2	841	527	gi 1196921	Unknown protein (Insertion sequence 18461)	90	70	315
48	27	20398	13757	gn F01024705	Lactate oxidase (Streptococcus faecalis)	90	80	1152
55	21	13977	18515	gn F01024705	Cpx protein (Bacillus subtilis)	90	75	1263
56	2	717	977	gi 1710133	Flagellar filament cap (Borrelia burgdorferi)	90	50	261
65	1	1	606	gi 1165303	L3 (Bacillus subtilis)	90	75	606
114	1	2	988	gi 153562	Aspartate beta-semialdehyde dehydrogenase (EC 1.2.1.11) (Streptococcus mutans)	90	80	987
120	1	1345	827	gi 407880	ORF1 (Streptococcus equisimilis)	90	75	519
159	12	7590	828	gi 143012	ORF synthetase (Bacillus subtilis)	90	84	609
166	4	4076	3282	gi 1661179	High affinity branched chain amino acid transport protein (Streptococcus mutans)	90	78	795
183	1	28	1395	gi 1048858	ATP-2,3-bisphosphotransferase (Lactococcus lactis)	90	76	1368
191	3	2891	1662	gi 149521	Cryptophan synthase beta subunit (Lactococcus lactis)	90	78	1230
198	2	1551	436	gi 2123242	(M004400) CcpA (Streptococcus mutans)	90	76	1116
305	1	37	783	gi 1573551	Asparagine synthetase A (anaA) (Haemophilus influenzae)	90	80	787
8	3	2285	3343	gi 144834	Purative (Lactococcus lactis)	89	78	1059
46	8	7577	7342	gi F045834A54	Ribosomal protein L19 - Bacillus stearothermophilus	89	76	216
49	9	8365	10142	gi 153722	RecP peptide (Streptococcus pneumoniae)	89	83	1380
51	14	18410	19447	gi 108857	ATP-D-fructose 6-phosphate 1-phosphotransferase (Lactococcus lactis)	89	81	1018
57	11	9686	10669	gn F01024705	ADP-forming MDH oxidase (Streptococcus mutans)	89	77	984
65	5	2418	2786	gi 1165307	L3 (Bacillus subtilis)	89	81	369
65	8	3806	4225	gi F04577R16	S05 RIBOSOMAL PROTEIN L16	89	82	420
65	18	8219	8719	gi 143417	Ribosomal protein S5 (Bacillus stearothermophilus)	89	76	501
73	9	6337	5315	gi 532204	Prs (Listeria monocytogenes)	89	70	1023
76	3	3360	1465	gn F01024705	LysA gene product (Bacillus subtilis)	89	76	1896
99	10	12818	11919	gi 153738	Membrane protein (Streptococcus mutans)	89	73	900
120	2	3552	3100	gi 407881	Arriquant capsid-like protein (Streptococcus equisimilis)	89	79	2753
122	5	4512	2791	gn F01024705	Unknown (Streptococcus pneumoniae)	89	81	1722

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

GenBank ID	Accession	Match gene name	% sim	% ident	length (nt)
176	1	5-oxopropyl-peptidase (Streptococcus pyogenes)	89	78	666
177	6	3050 3914 [gi 191423	89	71	865
181	8	4033 5751 [gi 149411	89	80	1719
211	4	3148 7293 [gi 532723	89	83	357
361	1	431 838 [gi 116922	89	70	408
34	17	11931 15535 [gi 20051916	88	78	1305
38	3	1446 14823 [gi 2028544	88	78	978
54	3	127 172 [gi 17003330	88	66	225
57	2	611 1468 [gi 170134943	88	75	858
65	13	5497 6089 [gi A231021808	88	75	573
65	20	9200 9500 [gi 2078381	88	83	471
78	3	3636 1108 [gi 170134943	88	80	2529
106	12	12965 12054 [gi 2407215	88	72	912
107	2	219 942 [gi 170134943	88	75	744
111	8	14073 10420 [gi 402363	88	74	3554
126	9	13086 12062 [gi 170134943	88	74	1035
140	17	19143 18874 [gi 1573659	88	61	270
144	1	354 555 [gi 170134943	88	75	162
148	4	2723 3493 [gi 1551672	88	65	771
160	8	5853 6278 [gi 1773267	88	65	426
177	4	1770 2885 [gi 149426	88	72	1116
211	6	4140 3613 [gi 532473	88	74	528
231	4	540 957 [gi 10146	88	78	378
260	5	2387 2998 [gi 116922	88	69	612
291	6	2047 3175 [gi 170134943	88	75	1359
319	4	658 317 [gi 603578	88	88	342
40	5	4353 4314 [gi 158972	87	56	162

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Config ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
49	10	10680	10929	gi 1196921	unknown protein Insertion sequence [S861]	87	72	270
65	7	3140	3808	gi 1165309	S3 [Bacillus subtilis]	87	73	669
65	15	6623	7039	gi 1044978	ribosomal protein S8 [Bacillus subtilis]	87	73	417
75	6	5411	6625	gi 1877422	galactokinase [Streptococcus mutans]	87	78	1215
80	2	703	2805	gi 19140116	elongation factor G [Bacillus subtilis]	87	76	2103
82	1	541	248	gi 1196921	unknown protein Insertion sequence [S861]	87	69	294
140	21	126033	23897	gi 19140116	elongation factor beta subunit [Bacillus subtilis]	87	74	1137
214	14	10441	8516	gi 2281305	glucose inhibited division protein homolog Gida [Lactococcus lactis cremoris]	87	75	1926
220	2	2742	874	gi 19140116	product highly similar to elongation factor EF-G [Bacillus subtilis]	87	73	1869
260	4	2096	2389	gi 1196921	unknown protein Insertion sequence [S861]	87	72	294
323	1	27	650	gi 1897795	30S ribosomal protein [Pediococcus acidilactici]	87	73	624
357	1	154	570	gi 1044978	ribosomal protein S8 [Bacillus subtilis]	87	73	417
49	11	10927	11445	gi 1196922	unknown protein Insertion sequence [S861]	86	69	519
59	12	7461	9274	gi 1951051	relaxase [Streptococcus pneumoniae]	86	68	1764
65	4	1553	2401	gi 1902759	ribosomal protein L2 [Bacillus stearothermophilus]	86	77	849
65	23	10957	11610	gi 144074	adenylate kinase [Lactococcus lactis]	86	76	654
82	4	4374	4856	gi 153745	.....	86	72	483
102	4	4270	4986	gi 19140116	.....	86	76	717
106	6	7824	6800	gi 19140116	.....	86	68	945
107	1	1	273	gi 19140116	.....	86	71	273
111	7	10437	8710	gi 19140116	.....	86	71	273
131	9	5704	4852	gi 1661193	.....	86	80	3733
134	7	6430	7980	gi 2308637	.....	86	71	813
146	11	7473	6583	gi 1591731	.....	86	73	1551
153	2	595	2010	gi 2180707	.....	86	72	851
154	1	2	1415	gi 1857246	.....	86	78	1416
					6-phosphogluconate dehydrogenase [Lactococcus lactis]	86	74	1434

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start	Stop	Match	Match accession	Match gene name	% sim	% ident	length (nt)
161	5	5025	6284	g1147529	Unknown [Streptococcus salivarius]		86	66	1260
184	1	2	1483	g1162467	HADP-dependent glyceraldehyde-3-phosphate dehydrogenase [Streptococcus sp.]		86	73	1482
210	8	3659	6571	g1151361	Translational initiation factor IF2 [Enterococcus faecium]		86	76	2513
250	1	2	187	g1157255	Asparagine synthetase A [aamA] [Haemophilus influenzae]		86	68	186
36	4	2644	3909	g11241909	Cell division protein [Enterococcus faecalis]		85	72	1266
38	4	2675	3587	g11204856	Putative ABC transporter subunit ComB [Streptococcus gordonii]		85	72	1113
38	5	3577	3915	g11205846	CoenC [Streptococcus gordonii]		85	80	339
57	5	2797	3789	g111974103170	IQG [Bacillus subtilis]		85	72	993
82	5	4915	6034	g11153746	Hamitol-phosphate dehydrogenase [Streptococcus mutans]		85	66	1140
81	15	11490	15793	g11143371	Phosphoribosyl aminimidazole synthetase [Pur-M] [Bacillus subtilis]		85	69	1104
87	2	1417	2388	g11184967	ScrR [Streptococcus mutans]		85	69	972
108	3	2668	3534	g11553586	OMP 13K protein [Enterococcus faecalis]		85	67	489
127	2	312	692	g11044989	Ribosomal protein S13 [Bacillus subtilis]		85	72	381
128	3	153	2409	g11665110	Tetrahydrofolate dehydrogenase/cyclohydrolase [Streptococcus thermophilus]		85	71	876
137	7	2962	4767	g111974100347	Nax -ATPase alpha subunit [Enterococcus hirae]		85	74	1806
170	2	2422	709	g111974102006	FUNCTION UNKNOWN, SIMILAR PRODUCT IN E. COLI, H. INFLUENZAE AND HAEMOPHILUS INFLUENZAE. [Bacillus subtilis]		85	70	1914
187	5	3760	4386	g11727436	Putative 20-Kb protein [Bacillus subtilis]		85	65	627
213	2	728	1873	g11163136	OMP-5 [Streptococcus pneumoniae]		85	67	1146
234	3	982	1255	g112293155	OMP8220 Y1A [Bacillus subtilis]		85	61	294
240	1	709	1307	g11435977	CTP synthetase [Bacillus subtilis]		85	70	1623
6	1	199	1521	g11508979	CTP-binding protein [Bacillus subtilis]		84	72	1333
10	4	4375	3843	g111974103982	Putative acylneuraminate lyase [Clostridium tertium]		84	70	933
14	1	63	2093	g11520753	DNA topoisomerase I [Bacillus subtilis]		84	69	2031
19	4	1793	2593	g11232484	(AF000508) RMase II [Bacillus subtilis]		84	68	801
20	17	17720	19687	g111974100584	Cell division protein [Bacillus subtilis]		84	71	1948
22	26	21721	22084	g11293163	Alanine dehydrogenase [Bacillus subtilis]		84	68	840



TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	Length (nt)
10	10	7730	6792	gi 100296	[fructokinase (Streptococcus pneumoniae)]			
31	9	5650	5300	gi 147194	[phnA protein (Escherichia coli)]	84	75	939
36	22	21551	20772	gi 1310631	[JVP binding protein (Streptococcus gordonii)]	84	71	351
48	4	2837	2505	gi 1882609	[6-phospho-beta-glucosidase (Escherichia coli)]	81	72	780
58	1	41	1516	gi 450849	[arabinose 5-phosphate isomerase (Streptococcus bovis)]	84	69	323
59	10	6715	7116	gi 1931053	[ORF10, putative (Streptococcus pneumoniae)]	84	73	1476
62	1	21	644	gi 1806487	[ORF23, putative (Lactococcus lactis)]	84	74	402
65	17	7779	8207	gi 1044980	[ribosomal protein L18 (Bacillus subtilis)]	84	66	624
65	21	9507	10397	gi 144073	[Secret protein (Lactococcus lactis)]	84	73	429
106	4	5474	2562	gi 193387	[carbamoyl-phosphate synthase (Lactobacillus plantarum)]	84	68	891
159	1	147	4	gi 1806487	[ORF21, putative (Lactococcus lactis)]	84	73	3213
163	4	4690	5310	gi 1293164	[AF008220] GSN synthase (Bacillus subtilis)]	84	63	144
192	1	46	1308	gi 1450556	[tripectidase (Lactococcus lactis)]	84	69	1221
348	1	671	6	gi 1787753	[AE002045] (346, 79 pct identical to 336 amino acids of ADH1_20000 SM: P20368 but has 10 additional N-ter residues (Escherichia coli)]	84	71	666
3	4	1572	3575	gi 1457466	[threop [BC 6.1.1.3] (Bacillus subtilis)]	83	65	2004
9	6	3833	3417	gi 121d100576	[single strand DNA binding protein (Bacillus subtilis)]	83	68	177
17	15	7426	8457	gi 150738	[coca protein (Streptococcus pneumoniae)]	83	66	1032
20	12	13860	14144	gi 121d100583	[unknown (Bacillus subtilis)]	83	61	285
23	4	3358	2806	gi 1788294	[AE002040] (235; This 235 aa orf is 40 pct identical (5 gaps) to 231 residues of an approx. 240 aa protein Y8C_2000 SM: P24237 (Escherichia coli)]	83	74	753
28	6	3104	3005	gi 1573455	[H. influenzae predicted coding region H0659 (Haemophilus influenzae)]	83	57	300
35	7	9108	3867	gi 3131707	[hypothetical nucleotide binding protein (Acetivibrio laticauda)]	83	63	1242
55	19	17932	17538	gi 1537085	[ORF_141 (Escherichia coli)]	83	59	405
55	20	18539	17919	gi 1494558	[orfX (Bacillus subtilis)]	83	69	621
65	6	2792	2142	gi 1165108	[L22 (Bacillus subtilis)]	83	64	348
68	6	6877	6683	gi 1213494	[immunoglobulin A1 protease (Streptococcus pneumoniae)]	83	54	195

TABLE 2

5. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	Length (nt)
87	15	15112	14771	gnl pid d121522	putative p202 protein [Bacillus subtilis]	83	54	342
86	12	6963	9631	gi 147394	5-oxopropyl-peptidase [Streptococcus pyogenes]	83	73	669
98	1	3	263	gi 1183885	glutamine-binding subunit [Bacillus subtilis]	83	55	281
120	4	7170	9233	gi 310620	zinc metalloprotease [Streptococcus gordonii]	83	72	1938
127	7	2998	4347	gi 1500567	M. Jannaschii predicted coding region M01465 [Methanococcus jannaschii]	83	72	1350
137	1	3	440	gi 1472918	γ-type Na-ATPase [Enterococcus hirae]	83	60	438
160	6	3466	4356	gi 1773265	ATPase, gamma-subunit [Streptococcus mutans]	83	67	891
214	4	2278	2964	gi 563279	transposase [Streptococcus pneumoniae]	83	72	687
226	3	2367	2020	gi 142154	chloroalcohol [Streptococcus pyogenes]	83	58	348
303	1	3	1049	gi 140046	phosphoglucose isomerase A (IA 1-46) [Bacillus stearothermophilus]	83	67	1047
6	17	15570	15318	gi 533147	glutamine synthetase [Bacillus subtilis]	82	64	1053
7	1	259	96	gi 141640	ribosomal protein L28 [Bacillus subtilis]	82	69	204
9	3	1479	1090	gi 385178	unknown [Bacillus subtilis]	82	46	390
9	7	4213	3899	gnl pid d100576	ribosomal protein S6 [Bacillus subtilis]	82	60	315
12	17	13422	14637	gnl pid d100571	unknown [Bacillus subtilis]	82	68	747
22	17	13422	14637	gnl pid d100571	putative [Bacillus subtilis]	82	68	747
22	18	14897	15658	gnl pid d101929	pyridine monophosphate kinase [Synecococcus sp.]	82	62	762
33	16	11471	10841	gnl pid d101190	ORF4 [Streptococcus mutans]	82	65	1416
35	9	7600	6255	gi 1881543	UDP-N-acetylglucosamine-2-epimerase [Streptococcus pneumoniae]	82	68	831
40	10	8003	7233	gi 1173519	riboflavin synthase beta subunit [Lactobacillus plantarum]	82	68	1146
48	32	23158	21437	gnl pid d10092	outer membrane protein [Campylobacter jejuni]	82	61	279
52	14	13833	14765	gi 142521	deoxyribidipyrrolidone photolase [Bacillus subtilis]	82	61	933
60	4	4737	1849	gnl pid d102221	AA0015100 uvaA [Deinococcus radiodurans]	82	66	2889
62	4	2131	1457	gi 2246748	AP009422 chloroalcohol reductase [Listeria monocytogenes]	82	63	975
71	11	18586	17818	gnl pid d122053	iso-1,4-α-D-glucan-1,4-glucosyltransferase [Streptococcus pneumoniae]	82	60	933
73	13	9222	7937	gnl pid d100586	unknown [Bacillus subtilis]	82	65	1386

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

TABLE 2

Contig ID	Start (bp)	Stop (bp)	Accession	Match gene name	% sim	% ident	Length (bp)
74	1	3771	[em P D d 01199]	alkaline amylopullulanase [Bacillus sp.]	82	68	3771
8	3	3494	[em P D d 00542]	alkaline protein product [Streptococcus thermophilus]	82	52	3494
86	11	10776	[gi 681583]	[5-ethylpyruvylshikimate-3-phosphate synthase [Lactococcus lactis]	82	67	1383
89	12	8495	[gi 40025]	homologous to E.coli 50K [Bacillus subtilis]	82	66	1428
115	9	10347	[em P D d 02090]	[AB001927] phospho-beta-galactosidase I [Lactobacillus gasseri]	82	74	1556
116	1	332	[em P D d 00595]	glyoxalase synthase [Bacillus subtilis]	82	71	1332
151	3	4657	[gi S06097 S060]	type I site-specific deoxyribonuclease (EC 3.1.21.3) CfrA chain S - Citrobacter freundii	82	66	1590
173	6	4183	[gi 2313836]	[AB000584] conserved hypothetical protein [Helicobacter pylori]	82	68	661
177	12	5481	[em P D d 01999]	[AB001341] NcrB [Escherichia coli]	82	58	1962
193	2	178	[gi S08564 P085]	ribosomal protein S9 - Bacillus stearothermophilus	82	70	359
245	2	258	[gi 146402]	EcoA type I restriction-modification enzyme S subunit [Escherichia coli]	82	68	568
9	5	3400	[em P D d 00576]	ribosomal protein S18 [Bacillus subtilis]	81	66	255
16	7	7484	[gi 1106074]	tryptophanyl-tRNA synthetase [Clostridium longispotum]	81	70	930
20	11	10398	[em P D d 00583]	transcription-repair coupling factor [Bacillus subtilis]	81	63	2515
38	2	1232	[gi 2058543]	putative DNA binding protein [Streptococcus gordonii]	81	63	375
45	2	3661	[gi 1460259]	isomerase [Bacillus subtilis]	81	67	1311
46	1	2	[gi 431231]	uracil permease [Bacillus caldolyticus]	81	61	1266
46	3	2455	[em P D d 00453]	mannosephosphate isomerase [Streptococcus mutans]	81	70	1014
54	2	1106	[gi 154752]	transport protein [Agrobacterium tumefaciens]	81	64	771
65	22	10306	[gi 144073]	[SecY protein [Lactococcus lactis]	81	66	516
89	4	3874	[gi 956886]	serine hydroxymethyltransferase [Bacillus subtilis]	81	69	1272
99	16	19126	[gi 2313526]	[AE000557] H. pylori predicted coding region HP0411 [Helicobacter pylori]	81	75	198
106	7	8373	[em P D d 019384]	[pyrK [Lactobacillus plantarum]	81	61	552
108	6	5054	[gi 1469939]	[group B oligonucleotide PspB [Streptococcus agalactiae]	81	66	1884
113	15	15889	[gi S09411 S094]	[ap0118 protein - Bacillus subtilis]	81	65	2385
128	5	3359	[gi 1665111]	[orf109] [Streptococcus thermophilus]	81	69	276

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
151	1	830	3211	gi 1046956	BooA type I restriction-modification enzyme R subunit [Escherichia coli]	81	59	2382
159	11	6722	7837	gi 2219788	GMP synthase [Bacillus subtilis]	81	69	1116
170	1	739	456	gnl P1D102906	FUNCTION UNKNOWN, [Bacillus subtilis]	81	55	282
191	2	1759	693	gi 149522	cryptophan synthase alpha subunit [Lactococcus lactis]	81	65	867
214	3	2290	1994	gi 157582	reverse transcriptase endonuclease [Drosophila virilis]	81	43	297
217	4	4415	4008	gi 466473	cellulose phosphorylase II [Bacillus stearothermophilus]	81	59	408
262	2	569	866	gi 153675	legatase 6-P kinase [Streptococcus mutans]	81	68	300
299	1	663	4	gnl P1D103194	3',5'-cyclic nucleotide phosphodiesterase [Salmonella enteritidis]	81	60	660
366	2	376	83	gi 149521	cryptophan synthase beta subunit [Lactococcus lactis]	81	65	294
12	10	8766	9242	gi 1216490	RNA/proteinase synthetase flavioprotein [Streptococcus mutans]	80	64	477
17	11	6050	5748	gnl P1D103162	unmated protein product [Streptococcus thermophilus]	80	67	303
17	16	8455	9046	gi 703126	hemagglutinin A receptor [Neisseria meningitidis]	80	59	612
18	3	2440	1613	gi 1531672	phosphate transport system ATP-binding protein [Methanococcus jannaschii]	80	58	828
27	3	4248	1579	gi 453109	galactose-4-epimerase [Bacillus subtilis]	80	69	2670
28	7	3671	1288	gi 1573660	H. influenzae predicted coding region HI0660 [Haemophilus influenzae]	80	63	364
32	2	902	1933	gnl P1D1026499	dihydroxyacetate dehydrogenase B [Lactococcus lactis]	80	66	1032
39	1	1	1266	gnl P1D1024078	hem [Lactococcus lactis]	80	63	1266
52	5	4261	3593	gi 1161864	ATP-binding subunit [Bacillus subtilis]	80	57	771
54	5	4250	4744	gi 2158820	(AF004225) Cux/CDB (18L); Cux/CDB homeoprotein [Mus musculus]	80	60	195
59	11	7109	7486	gi 1951052	ORF5, putative [Streptococcus pneumoniae]	80	68	378
65	3	1230	1550	gi A028154305	ribosomal protein L23 - Bacillus stearothermophilus	80	69	321
65	12	5174	5903	gi A028194305	ribosomal protein L24 - Bacillus stearothermophilus	80	70	330
66	9	9484	10687	gi 2313816	(AE000584) conserved hypothetical protein [Haemophilus pylori]	80	66	804
82	2	648	2438	gi 222991	hemolysin transport protein [Bacillus stearothermophilus]	80	65	1791
85	1	950	630	gi 528995	polyketide synthase [Bacillus subtilis]	80	46	321
89	8	6670	5779	gi 853776	peptide chain release factor 1 [Bacillus subtilis]	80	63	1092
93	12	8718	7438	gnl P1D101959	hypothetical protein [Synedococcus sp.]	80	60	1281

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig	ORF	Start	Stop	match accession	match gene name	% sim	% ident	length (nt)
106	5	6854	5751	gnl PID e19386	glutaminase of carbamoyl-phosphate synthase (Lactobacillus plantarum)	80	65	1104
109	2	2140	1450	gnl 40056	[por] gene product (Bacillus subtilis)	80	59	711
124	9	4246	3953	gnl PID d02254	30S ribosomal protein S16 (Bacillus subtilis)	80	65	294
128	8	5148	6428	gnl 2281308	phosphoenolcarboxylate (Lactococcus lactis cremoris)	80	66	1281
137	139	12665	11376	gnl 159109	ADP-dependent glutamate dehydrogenase (Gardinia intestinalis)	80	68	1290
140	19	18659	19457	gnl 517210	putative transposase (Streptococcus pyogenes)	80	70	243
158	2	2474	984	gnl 1877423	galactose-1-P-uridylyl transferase (Streptococcus mitis)	80	65	1491
171	110	7174	7728	gnl 379860	cytoplasmic C-associated protein (Mus musculus)	80	60	255
181	1	2	619	gnl 149395	JAC (Lactococcus lactis)	80	66	618
313	1	27	539	gnl 143467	ribosomal protein S4 (Bacillus subtilis)	80	70	513
329	2	1652	858	gnl 533080	BeP protein (Streptococcus pyogenes)	80	63	795
371	1	2	958	gnl 422360	ClpC adenosine triphosphatase (Bacillus subtilis)	80	58	957
8	7	4312	5580	gnl 149435	putative (Lactococcus lactis)	79	64	1269
23	1	1175	135	gnl 1542975	JACB (Thermotogaobacterium thermophilum)	79	61	1041
33	14	9244	8201	gnl PID e253831	UOP-glucose 4-epimerase (Bacillus subtilis)	79	62	1044
36	3	1242	2633	gnl PID e24218	[fzA] Enterococcus faecalis	79	58	1392
38	13	7155	8378	gnl 405134	secreted kinase (Bacillus subtilis)	79	58	1224
55	7	9011	8229	gnl 1146234	dihydrodipicolinate reductase (Bacillus subtilis)	79	56	783
65	19	4661	8915	gnl 2078380	ribosomal protein L30 (Staphylococcus aureus)	79	68	255
69	4	3678	2128	gnl PID e11452	unknown (Bacillus subtilis)	79	64	1351
69	9	7881	7279	gnl 677890	hypothetical protein (Staphylococcus aureus)	79	59	603
72	10	4491	3783	gnl PID d101091	hypothetical protein (Synecocystis sp.)	79	62	1293
80	3	2906	7100	gnl 141342	polymerase III (Bacillus subtilis)	79	65	4395
82	14	13326	15689	gnl PID e255093	hypothetical protein (Bacillus subtilis)	79	65	2364
86	13	12233	11118	gnl 1883582	prephenate dehydrogenase (Lactococcus lactis)	79	58	1116
92	3	940	1734	gnl 537286	triosphosphate isomerase (Lactococcus lactis)	79	65	795
98	6	4023	4742	gnl PID d100262	lipo protein (Salmonella typhimurium)	79	63	720

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
99	12	16315	14150	gi 153736	ecolactide (Streptococcus mitis)	79	64	2166
107	7	5684	6406	gi 1460080	D-alanine-D-alanine ligase-related protein (Enterococcus faecalis)	79	58	723
113	9	6858	8303	gi 1465882	ppp1, B1496_C2,189 Mycolactium leprae	79	64	1446
151	10	13424	12213	gi 4506666	3-phosphoglycerate kinase (Thermotoga maritima)	79	60	1212
162	2	1158	3017	gi 505700	CcpA (Staphylococcus aureus)	79	67	1860
177	5	2876	3052	gi 912423	putative (Lactococcus lactis)	79	61	177
187	8	4198	4503	gi 147404	putative (Lactococcus lactis)	79	61	366
189	3	2728	2907	gm PF01614149	putative ATP-binding protein of ABC-type (Bacillus subtilis)	79	53	180
191	5	4249	3449	gi 149519	indoleglycerol phosphate synthase (Lactococcus lactis)	79	61	762
211	3	3805	2737	gi 147404	mannose permease subunit II-M-Han (Escherichia coli)	79	66	801
212	3	3863	3621	gm PF04205004	glutaredoxin-like protein (Lactococcus lactis)	79	57	933
215	1	987	715	gi 2253242	[AF008220] arginine succinate synthase (Bacillus subtilis)	79	56	243
323	2	530	781	gi 897795	30S ribosomal protein (Pediococcus acidilactici)	79	64	273
380	1	694	2	gi 1146680	[polynucleotide phosphorylase (Bacillus subtilis)]	79	67	252
384	2	655	239	gi 143328	JobP protein (put.: putative (Bacillus subtilis))	79	64	653
6	3	2820	1091	gi 852767	[UDP-N-acetylglucosamine 1-carboxyvinyltransferase (Bacillus subtilis)]	78	62	1272
8	1	50	1786	gi 145432	putative (Lactococcus lactis)	78	63	1737
9	1	351	124	gi 897793	y98 gene product (Pediococcus acidilactici)	78	59	228
15	8	7364	8314	gm PF0105085	lysatein synthetase A (Bacillus subtilis)	78	63	951
20	10	9738	10310	gm PF0105083	stage V sporulation (Bacillus subtilis)	78	58	573
20	16	17765	17713	gi 149105	[polynucleotide phosphorylase (Bacillus subtilis)]	78	59	549
22	22	17388	18416	gm PF0101315	YqfE (Bacillus subtilis)	78	60	1029
22	37	20571	20612	gi 139916	alanine dehydrogenase (Bacillus subtilis)	78	59	360
34	8	7407	7105	gi 11015	aspartate-tRNA ligase (Escherichia coli)	78	55	302
35	8	6247	5196	gi 1457644	CapB (Staphylococcus aureus)	78	60	1062

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	Length (nt)
40	11	9287	8001	gi 1173518	GTP cyclohydrolase II/ 3,4-dihydroxy-2-butanone-4-phosphate synthase (Actinobacillus pleuropneumoniae)	78	56	1287
48	31	12422	23183	gi 2314330	(AD006623) glutamine ABC transporter, ATP-binding protein (glnQ) (Haemophilus pylori)	78	58	762
52	2	2101	1430	gi 1183887	integral membrane protein [Bacillus subtilis]	78	54	672
55	14	13405	12707	gi 1400028	(AB02150) Thp [Bacillus subtilis]	78	56	894
55	17	16637	15612	gi 17031027	hypothetical protein [Bacillus subtilis]	78	51	1026
71	14	19756	19598	gi 17967	actin channel alpha-10 subunit (Homo sapiens)	78	57	159
74	11	15031	14018	gi 1573279	Holliday junction DNA helicase (rvb8) (Haemophilus influenzae)	78	57	1014
75	9	6623	7877	gi 1577423	galactose-1-P-uridylyl transferase [Streptococcus mutans]	78	62	1350
81	12	12125	11906	gi 1573607	L-fucose isomerase (fucI) [Haemophilus influenzae]	78	66	1282
82	5	2037	4417	gi 153744	ORF X, putative [Streptococcus mutans]	78	64	1995
83	18	16926	18500	gi 143373	phosphoribosyl aminimidazole carboxy formyl formyltransferase/inosine monophosphate cyclohydrolase (Pur-H10) [Bacillus subtilis]	78	63	1575
83	20	20212	20775	gi 143364	phosphoribosyl aminimidazole carboxy formyl transferase (Pur-H8) [Bacillus subtilis]	78	64	564
92	2	165	878	gi 1401190	ORF2 [Streptococcus mutans]	78	62	714
98	8	5863	6902	gi 2311287	(AF031188) repressor 2 [Bacillus subtilis]	78	63	1047
113	3	1071	2741	gi 580914	dnaX [Bacillus subtilis]	78	64	1671
127	4	1133	2071	gi 142165	DNA polymerase alpha-core-subunit [Bacillus subtilis]	78	59	939
132	1	2782	497	gi 1561763	pullulanase (Bacteroides thetaiotaomicron)	78	58	2286
135	4	2658	2537	gi 1788036	(AB00289) NH <sub>2</sub> -dependent MD synthetase [Escherichia coli]	78	66	840
140	24	26853	25423	gi 1100077	phospho-beta-glucosidase (Clostridium longosporum)	78	64	1431
150	5	4690	4516	gi 149464	amino peptidase [Lactococcus lactis]	78	42	177
152	1	1	795	gi 1639915	NADH dehydrogenase subunit [Rhizogonia alata]	78	55	795
162	4	4997	4110	gi 140123528	putative Thp protein [Bacillus subtilis]	78	64	888
181	10	8651	7947	gi 149402	lactose repressor (lact <sup>+</sup> alt.) [Lactococcus lactis]	78	48	705
200	4	3627	4598	gi 140100172	invertase [Symonias mobilis]	78	61	1332
203	3	3230	3015	gi 1174237	CysK [Pseudomonas fluorescens]	78	57	216

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
210	9	6789	7172	gi 1569902	ORF6 gene product [Bacillus subtilis]	78	42	384
214	6	3810	2797	gn J010249	P <sub>1</sub> hemolysin o-nitrolysoprotein endopeptidase; P36175 (660) [Bacillus subtilis]	78	60	1014
214	13	6322	8163	gi 1377831	unknown [Bacillus subtilis]	78	62	1842
217	1	9	2717	gi 1488430	alcohol dehydrogenase 2 [Bacillus histolytica]	78	64	2709
222	3	2316	3098	gi 1571047	spore germination and vegetative growth protein (gacC2) [Haemophilus influenzae]	78	65	783
268	1	742	8	gi 517210	putative transposase [Streptococcus pyogenes]	78	65	735
276	1	233	793	gn J0100306	ribosomal protein L1 [Bacillus subtilis]	78	65	531
312	3	1567	1079	gi 287261	conG ORF2 [Bacillus subtilis]	78	54	459
339	1	117	794	gi 1916729	cadD [Staphylococcus aureus]	78	53	678
342	2	742	265	gi 1802439	phosphatidylglycerophosphate synthase [Bacillus subtilis]	78	59	498
383	1	737	3	gi 1146880	polynucleotide phosphorylase [Bacillus subtilis]	78	64	735
7	15	11018	11018	gi 1319855	carboxyltransferase beta subunit [Streptococcus RCC943]	77	63	906
8	2	1698	2255	gi 149433	putative lactococcus lactis]	77	59	558
17	134	6948	7250	gi 5207218	conA protein [Streptococcus pneumoniae]	77	60	603
30	12	9761	8967	gi 1000451	TrpP [Bacillus subtilis]	77	43	795
36	14	11231	11231	gi 1537786	phosphoglyceronate (gapA) [Haemophilus influenzae]	77	64	711
55	3	3836	6096	gi 1708640	YaaB [Bacillus subtilis]	77	55	261
61	8	8377	8054	gi 1890649	multidrug resistance protein LmrA [Lactococcus lactis]	77	51	324
65	2	607	1254	gi 140103	ribosomal protein L4 [Bacillus stearothermophilus]	77	63	648
68	8	7509	7240	gi 47551	IMP [Streptococcus suis]	77	68	270
69	1	1083	118	gn J011493	unknown [Bacillus subtilis]	77	57	988
77	5	4593	4026	gn J0215178	hypothetical 12.2 kd protein [Bacillus subtilis]	77	60	558
83	14	13104	14552	gi 1509947	amidophosphoryltransferase [Bacillus subtilis]	77	77	1449
94	4	3006	5444	gn J033895	(J0300496) cyclic nucleotide-gated channel beta subunit [Rattus norvegicus]	77	66	2439
96	11	8518	8880	gi 1511879	ORP 1 [Lactococcus lactis]	77	62	363
99	11	14082	12799	gi 153737	leuarg-binding protein [Streptococcus mutans]	77	61	1284





TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	length (nt)
22	37	27756	28171	gn1 P1D e13189	translation initiation factor IF3 (AA 1-172) [Bacillus stearothermophilus]	76	61	416
35	6	1369	2682	gi11773346	CapsG [Staphylococcus aureus]	76	61	1188
48	28	21113	21787	gi12314328	(AE000623) glutamine ABC transporter, permease protein (glnP) [Helicobacter pylori]	76	52	675
52	12	13281	13786	gi1432321	acylribitolpyrimidine photocyclase [Bacillus subtilis]	76	58	506
55	110	11131	10571	gn1 P1D e183110	femp [Staphylococcus aureus]	76	61	551
57	8	7824	4559	gi1250561	o188 [Escherichia coli]	76	61	321
62	5	1406	2095	gn1 P1D e113024	hypothetical protein [Bacillus subtilis]	76	59	712
65	9	1232	44	gi140146	129 protein (AA 1-66) [Bacillus subtilis]	76	58	219
68	2	1328	2371	gn1 P1D e184233	anabolic ornithine carbamoyltransferase [Lactobacillus plantarum]	76	61	1044
69	8	7297	8005	gn1 P1D d101420	pyrimidine nucleoside phosphorylase [Bacillus stearothermophilus]	76	61	1293
73	12	7839	7267	gn1 P1D e143629	unknown [Mycobacterium tuberculosis]	76	53	573
74	5	8433	7039	gn1 P1D d102048	C. thermocellus beta-glucosidase; P26208 (965) [Bacillus subtilis]	76	60	1395
80	5	7643	7936	gi12314030	(AE000599) conserved hypothetical protein [Helicobacter pylori]	76	61	294
82	15	14019	16936	gi11573900	D-alanine permease (dagh) [Haemophilus influenzae]	76	56	978
83	19	18616	19884	gi1143374	phosphoribosyl glycinaamide synthetase (PUB-D, 902 start codon) [Bacillus subtilis]	76	60	1269
86	14	13409	13231	gi1143366	AcP [Bacillus subtilis]	76	58	1179
87	1	3	1442	gi11518604	sucrose-6-phosphate hydrolase [Streptococcus mutans]	76	59	1440
87	16	15754	1510	gn1 P1D e22500	putative oak protein [Bacillus subtilis]	76	56	645
93	4	1769	1539	gi11574820	1,4-alpha-glucan branching enzyme (glgB) [Haemophilus influenzae]	76	46	231
94	3	51	385	gi1144313	6.0 kd ORF (P1amid ColE1)	76	73	315
116	2	2151	1678	gi1153841	pneumococcal surface protein A [Streptococcus pneumoniae]	76	59	474
123	6	3442	3875	gi11314297	ClpC ATPase [Listeria monocytogenes]	76	59	2454
126	2	2156	2932	gn1 P1D d101328	g102 [Bacillus subtilis]	76	61	777
128	10	6973	7797	gi1944944	purine nucleoside phosphorylase [Bacillus subtilis]	76	60	825
131	11	6186	5812	gi11674310	(AE000058) Mycoplasma pneumoniae, 60205 homolog, from M. genitalium [Mycoplasma pneumoniae]	76	47	375

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig	ORF	Start (nt)	Stop (nt)	Accession	match	match gene name	% sim	% Ident	length (nt)
139	4	3641	3192	gi 2293302	100	YtaA [Bacillus subtilis]	76	53	450
140	114	14872	12236	gi 1184680	100	polynucleotide phosphorylase [Bacillus subtilis]	76	62	2337
143	2	2583	3305	gi 143795	100	transfer RNA-Tyr synthetase [Bacillus subtilis]	76	61	1221
170	6	5095	6114	gn P100559	100	YggG [Bacillus subtilis]	76	44	1020
180	2	1927	557	gi 40019	100	ORF 821 [aa 1-821] [Bacillus subtilis]	76	53	1371
191	7	5815	5228	gi 551880	100	anthranilate synthase beta subunit [Lactococcus lactis]	76	61	588
195	3	3829	2444	gi 2149905	100	D-glutamic acid adding enzyme [Enterococcus faecalis]	76	60	1186
200	3	1934	3627	gi 431272	100	lysin protein [Bacillus subtilis]	76	58	1716
201	1	431	207	gi 2308998	100	desferal glucosidase hnsS [Streptococcus sulci]	76	57	225
214	2	1283	2380	gi 643278	100	transposase [Streptococcus pneumoniae]	76	55	1098
225	3	2138	3411	gi 1552775	100	ATP-binding protein [Escherichia coli]	76	56	1074
233	1	2	724	gi 1163115	100	neuraminidase B [Streptococcus pneumoniae]	76	60	723
347	1	523	38	gi 572033	100	ORF 036 [Escherichia coli]	76	60	486
356	2	842	165	gi 2149905	100	D-glutamic acid adding enzyme [Enterococcus faecalis]	76	61	678
366	3	734	346	gi 1189520	100	phosphatidylserine synthase [Lactococcus lactis]	76	69	387
5	8	12599	11484	gi 1574293	100	fibrillar transcription regulation repressor (p18) [Haemophilus influenzae]	75	61	1116
6	13	12553	11894	gn P100550	100	YggH [Bacillus subtilis]	75	51	660
9	10	7282	4062	gi 162538	100	aspartate aminotransferase [Bacillus sp.]	75	55	1221
10	12	8080	7940	gi 109493	100	SCRP methylase [Lactococcus lactis]	75	56	141
18	5	4266	3301	gn P100519	100	YggH [Bacillus subtilis]	75	52	946
22	4	1818	2728	gi 173157	100	orf-X; hypothetical protein; Method: conceptual translation supplied by [Bacillus subtilis]	75	62	891
30	111	9015	7828	gi 153801	100	enzyme ser-II [Streptococcus mutans]	75	64	1188
31	5	2782	2030	gi 2293311	100	ORF 820 putative thioredoxin [Bacillus subtilis]	75	53	1338
32	9	7484	8359	gn P100560	100	formamidopyrimidine-DNA glycosylase [Streptococcus mutans]	75	61	876
33	4	1725	1448	gi 143276	100	ipa-32r gene product [Bacillus subtilis]	75	53	288
33	10	4470	5769	gi 531105	100	unknown [Bacillus subtilis]	75	56	702

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
31	12	8878	7183	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	56	306
36	1	181	2	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	43	180
38	12	14510	15179	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	56	870
48	33	23398	24066	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	55	669
51	1	2	319	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	55	318
51	10	8318	11683	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	50	3366
54	18	19566	20759	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	58	1194
57	9	8448	7822	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	50	627
65	14	4072	6356	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	64	285
70	4	3071	2472	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	57	600
71	24	30399	29404	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	57	996
73	2	910	455	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	57	456
79	1	1810	491	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	59	1320
82	6	6360	6536	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	55	177
83	6	1938	2975	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	56	1038
93	11	7368	5317	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	58	2052
93	13	9409	8599	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	54	711
95	1	1795	47	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	57	1749
103	2	362	1186	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	64	825
104	1	691	915	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	54	225
113	5	2951	3893	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	55	933
121	1	320	1390	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	58	1071
127	6	2614	2000	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	44	387
137	18	10082	10687	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	32	608
149	11	8499	9318	gi 1208739	hypothetical protein (Campylobacter jejuni)	75	55	840

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	Start (nt)	Stop (nt)	match position	match gene name	% sim	% ident	length (nt)
151	6	9100	g1146467	HsdS polypeptide, part of CtrA family [Citrobacter freundii]	75	57	1428
158	3	986	g1146467	UMP-glucose 4-epimerase [Bacillus subtilis]	75	63	984
172	8	5653	g1142978	glycerol dehydrogenase [Bacillus stearotherophilus]	75	56	1122
172	9	7119	g1142978	glycerol dehydrogenase [Bacillus stearotherophilus]	75	58	2592
173	1	361	g1142978	unknown [Mycobacterium tuberculosis]	75	50	183
185	3	3066	g1142978	apennidine/potracine transport ATP-binding protein (pota) [Memophilus	75	56	1053
191	6	4235	g1149518	phosphoribosyl anthranilate transferase [Lactococcus lactis]	75	6	1023
226	2	1774	g1142978	homolog of E. coli ribosomal protein L21 [Bacillus subtilis]	75	65	594
231	1	153	g1140173	homolog of E. coli ribosomal protein L21 [Bacillus subtilis]	75	57	153
234	1	2	g112293259	unknown protein [Bacillus subtilis]	75	59	417
279	1	552	g11118198	unknown protein [Bacillus subtilis]	75	50	402
281	7	2558	g1140011	ORF1 (AA 1-161) [Bacillus subtilis]	75	58	492
375	2	137	g1140137	ORF13 [Bacillus subtilis]	74	53	840
6	10	10721	g112293259	ORF13 [Bacillus subtilis]	74	60	1371
7	6	4682	g11354211	ORF13-like protein [Bacillus subtilis]	74	54	915
18	4	3341	g112293259	glutamate-aminopeptidase [Lactococcus lactis]	74	59	1086
21	6	5885	g11072381	glutamate-aminopeptidase [Lactococcus lactis]	74	46	192
24	2	739	g112314762	glutamate-aminopeptidase [Lactococcus lactis]	74	63	366
25	1	2	g112314762	glutamate-aminopeptidase [Lactococcus lactis]	74	57	1533
38	18	11432	g11537034	ORF_0488 [Escherichia coli]	74	53	2256
48	10	8924	g11531069	ORF_0488 [Escherichia coli]	74	53	2256
55	11	11964	g11401	ORF_0488 [Escherichia coli]	74	53	2256
61	2	1792	g112293259	ORF_0488 [Escherichia coli]	74	53	2256
76	10	5414	g11401	ORF_0488 [Escherichia coli]	74	53	2256
83	2	666	g11401	ORF_0488 [Escherichia coli]	74	53	2256
86	9	8985	g11682585	ORF_0488 [Escherichia coli]	74	53	2256

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
102	5	5065	gi 141394	GMP-PRPP transferase [Bacillus subtilis]	74	57	648
103	5	4346	gi 141394	GMP-PRPP transferase [Bacillus subtilis]	74	62	1098
108	7	6884	gi 141394	Yton protein [Bacillus subtilis]	74	56	729
111	2	478	gi 141394	Methyltransferase [Bacillus subtilis]	74	45	333
133	2	1340	gi 141394	YogZ [Bacillus subtilis]	74	60	462
137	9	6167	gi 141394	Hypothetical protein [Bacillus subtilis]	74	53	631
149	4	3098	gi 141394	Ne <sup>+</sup> ATPase subunit D [Enterococcus faecalis]	74	55	876
157	2	243	gi 1573646	High level aminoglycoside resistance [Bacillus subtilis]	74	48	582
164	6	3515	gi 141031	Imethylated-DNA-protein-cysteine methyltransferase (dat1) [Haemophilus influenzae]	74	48	735
167	7	5446	gi 1413927	lpa-3r gene product [Bacillus subtilis]	74	55	246
171	1	1818	gi 14102251	Beta-galactosidase [Bacillus clausii]	74	62	1818
172	4	1064	gi 1466474	Cellulose phosphorylase enzyme [Bacillus stearothermophilus]	74	50	1325
185	1	326	gi 1573646	Hg(2+) transport ATPase protein C (mgC) (SP-P22037) [Haemophilus influenzae]	74	68	324
188	2	1089	gi 1573608	ATP dependent translocator homolog (msbA) [Haemophilus influenzae]	74	44	930
189	11	6491	gi 1661199	Sakacin A production response regulator [Streptococcus mutans]	74	60	684
210	2	250	gi 1229207	(AP008220) YnaQ [Bacillus subtilis]	74	60	768
261	1	836	gi 1666983	Putative ATP binding subunit [Bacillus subtilis]	74	55	645
263	3	1619	gi 1663232	Similarity with S. cerevisiae hypothetical 137.7 kb protein in subtelomeric repeat region [Saccharomyces cerevisiae]	74	42	2037
265	2	844	gi 149272	Asparaginase [Bacillus licheniformis]	74	64	384
368	1	1	gi 1602298	unknown [Saccharomyces cerevisiae]	74	39	942
7	16	11357	gi 14101324	YcaH [Bacillus subtilis]	73	57	1437
17	10	5706	gi 14101324	Unnamed protein product [Streptococcus thermophilus]	73	47	258
31	2	532	gi 14101324	Single strand RNA binding protein [Bacillus subtilis]	73	55	279
32	6	566	gi 14101315	YqgC [Bacillus subtilis]	73	58	528
34	15	10281	gi 14102151	(AB001684) OmpC [Chlorella vulgaris]	73	46	492

TABLE 2

5. pneumoniae Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	match	match gene name	% sim	% ident	length (nt)
40	12	8746	9236	gi1177351		ribosomal protein alpha subunit [Actinobacillus pleuropneumoniae]	73	55	651
55	2	3592	839	gn1 P D d 01887		cation-transporting ATPase PacL [Synecocystis sp.1]	73	60	2784
55	18	11494	15886	gn1 P D e 65580		unknown [Mycobacterium tuberculosis]	73	52	909
65	16	7213	7767	gi1103419		ribosomal protein L6 [Bacillus stearothermophilus]	73	60	555
68	3	3300	3659	gn1 P D e 65983		haer [Lactobacillus casei]	73	52	360
70	10	5557	6195	gn1 P D e 65761		envelope protein (human immunodeficiency virus type 1)	73	60	177
71	4	8133	8262	gn1 P D e 22063		ss-1,4-galactosyltransferase [Streptococcus pneumoniae]	73	45	2130
72	1	3	851	gi12293177		AP0802201 transporter [Bacillus subtilis]	73	50	849
76	7	7019	6195	gn1 P D d 01325		YqjF [Bacillus subtilis]	73	66	825
76	12	10009	9533	gi1573086		uridine kinase (uridine monophosphokinase) (uak) [Mycobacterium tuberculosis]	73	54	477
80	7	3113	3372	gi13377823		aminopeptidase [Bacillus subtilis]	73	60	1260
97	5	3389	1668	gn1 P D d 01954		dihydroxyacid dehydratase [Synecocystis sp.]	73	54	1722
98	9	6912	7619	gn1 P D e 14991		PrE [Mycobacterium tuberculosis]	73	54	708
108	11	10928	10440	gn1398109		regulatory protein [Enterococcus faecalis]	73	63	591
128	6	3632	4222	gi11685111		orf1091 [Streptococcus thermophilus]	73	54	489
138	2	1575	394	gi1147326		transport protein [Escherichia coli]	73	60	1182
140	13	12538	11903	gn1 E S3402 E534		serine O-acetyltransferase (EC 2.3.1.30) - Bacillus stearothermophilus	73	55	636
142	5	5700	4991	gn1 P D e 23251		putative 30S protein [Bacillus subtilis]	73	50	711
164	4	2323	2790	gi1592076		hypothetical protein (SP-P25763) [Mycobacterium tuberculosis]	73	52	468
164	8	6815	5566	gi1410137		ORF13 [Bacillus subtilis]	73	56	732
170	5	4394	5302	gn1 P D d 00955		homologue of unidentified protein of E. coli [Bacillus subtilis]	73	46	909
178	7	3890	4855	gi146242		modulation protein B, 5' and 3' noncoding region	73	56	963
204	6	5096	4278	gn1 P D e 214719		PLGA protein [Bacillus thuringiensis]	73	43	819
213	2	832	2037	gi1565296		ribosomal protein S1 homologue; sequence specific DNA-binding protein	73	55	1206
211	2	84	287	gi140173		homolog of E. coli ribosomal protein L21 [Bacillus subtilis]	73	61	204
217	1	2	505	gi1177351		adenine phosphoribosyltransferase [Escherichia coli]	73	51	504

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

TABLE 2

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
269	1	2	691	gnl pfid d101328	YqjX [Bacillus subtilis]	73	36	690
289	2	1272	832	pir A02771 P7MC	ribosomal protein L7/L12 - Micrococcus luteus	73	66	441
343	1	14	484	gll 1788125	(A0800276) hypothetical 30.4 kD protein in mmd-capC intergenic region [Escherichia coli]	73	47	471
358	1	222	4	g12149905	D-glutamic acid adding enzyme [Enterococcus faecalis]	73	50	219
7	5	3165	4691	gnl pfid d101833	amidase [Synecococcus sp.]	72	50	1527
7	7	7195	7647	gll 146976	nusB [Escherichia coli]	72	54	453
7	17	11743	13300	gnl pfid e289141	similar to hydromyristoyl-acyl carrier protein dehydratase [Bacillus subtilis]	72	55	444
22	19	15637	16274	gnl pfid d101925	ribosome releasing factor [Synecococcus sp.]	72	51	588
33	17	112111	11425	gnl pfid d101190	ORF3 [Streptococcus mutans]	72	55	687
34	7	7147	7672	gll 1394501	ribosomal protein S10 [Thermus thermophilus]	72	52	1521
38	23	15372	16085	pir H64108 H641	L-ribulose-phosphate 4-epimerase (arid) homolog - Haemophilus influenzae (strain Rd RM20)	72	54	714
39	5	5094	6905	gnl pfid e254877	unknown [Mycobacterium tuberculosis]	72	36	1812
40	6	4469	4636	gll 153672	lactose repressor [Streptococcus mutans]	72	58	168
48	2	1459	1253	gll 310380	Inhibin beta-A-subunit [Pisus arvensis]	72	33	207
48	29	21729	22474	gll 2314235	(A0000623) glutamine ABC transporter, permease protein (gluP) [Helicobacter pylori]	72	49	636
50	5	4529	3288	gll 1750108	YnhA [Bacillus subtilis]	72	50	1242
51	3	1044	2282	gll 2293230	(A0008220) YtkB [Bacillus subtilis]	72	54	1239
52	13	11681	11938	gll 142521	deoxyribodipyrimidine photolyase [Bacillus subtilis]	72	45	258
55	1	841	35	gll 882518	ORF_0304; ORF start [Escherichia coli]	72	59	807
75	5	2832	3191	gnl pfid e205886	mercuric resistance operon regulatory protein [Bacillus subtilis]	72	44	360
76	6	6229	5771	gll 142450	ahcC protein [Bacillus subtilis]	72	53	459
79	5	5065	4592	gll 2293279	(A0008220) YteG [Bacillus subtilis]	72	46	474
87	14	14726	13109	gnl pfid e23502	putative prfA protein [Bacillus subtilis]	72	52	2418
91	1	444	662	gll 500891	YK001 [Candida glabrata]	72	50	219
91	7	4516	4764	gll 828615	skeletal muscle sodium channel alpha-subunit [Equus caballus]	72	38	249



TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start	Stop	Match accession	Match gene name	% sim	% Ident	Length (nt)
95	2	2064	1717	gn P1D e323527	putative Asp2 protein [Bacillus subtilis]	72	40	288
109	1	1462	118	gi 143331	alkaline phosphatase regulatory protein [Bacillus subtilis]	72	52	1335
126	1	3	2192	gn P1D e101831	glutamine-binding periplasmic protein [Synechocystis sp.]	72	46	2190
130	3	1735	2478	gn P1D e215396	(AP015755) carboxypeptidase [Bacillus subtilis]	72	53	744
137	6	2585	2929	gi 472922	v-type Na-ATPase [Enterococcus faecalis]	72	48	399
140	10	3601	9201	gi 43224	URF 4 [Synechococcus sp.]	72	45	600
146	5	1596	1247	gn P1D e324945	hypothetical protein [Bacillus subtilis]	72	56	1002
147	2	2084	1083	gn P1D e325016	hypothetical protein [Bacillus subtilis]	72	56	1002
147	5	6156	5146	gi 472327	TPP-dependent actoin dehydrogenase beta-subunit [Clostridium magnum]	72	56	1011
148	8	2381	6433	gi 574332	IND(PH-dependent dihydroxyacetone-phosphate reductase [Bacillus subtilis]	72	54	1053
148	14	10256	9675	gn P1D e101319	Yogw [Bacillus subtilis]	72	50	582
159	8	4005	4949	gi 1788770	(AE000330) o463; 24 pct identical (44 gaps) to 338 residues from penicillin-binding protein 4*, PBP4-BACDU SH: P22959 (451 aa) [Escherichia coli]	72	43	945
172	10	9907	10620	gi 703387	unknown [Saccharomyces cerevisiae]	72	55	716
220	3	2862	3602	gi 1574175	hypothetical [Haemophilus influenzae]	72	50	741
267	1	3	449	gi 290513	[470 [Escherichia coli]	72	48	447
281	2	895	510	gn P1D e100964	homologue of acetylcholinesterase 2 alpha and beta subunits <i>lysc</i> of <i>B. subtilis</i> [Bacillus subtilis]	72	45	360
290	1	1018	14	gi 474195	This ORF is homologous to a 40.0 kd hypothetical protein in the htr-8 3' region from <i>E. coli</i> , Accession Number X61000 [Mycoplasma-like organism]	72	54	1005
300	1	63	587	gi 746399	transcription elongation factor [Escherichia coli]	72	50	525
316	1	1726	4	gn 158127	protein kinase C [Drosophila melanogaster]	72	40	1323
342	1	227	3	gn P1D e101164	unknown [Bacillus subtilis]	72	54	229
354	1	1	1005	gn P1D e102048	C. thermocellus beta-glucosidase; P22028 [985] [Bacillus subtilis]	72	52	1005
6	10	8134	10457	gn P1D e264229	unknown [Mycobacterium tuberculosis]	71	57	2314
7	20	16231	15444	gn 180046	3-oxoacyl-(acyl-carrier protein) reductase [Cuphea lanceolata]	71	52	768
15	1	1297	2	gn P1D e100571	replicative DNA helicase [Bacillus subtilis]	71	51	1299
15	4	4435	3869	gi 499384	orf189 [Bacillus subtilis]	71	47	567

TABLE 2

S. pneumoniae - Native coding regions of novel proteins similar to known proteins

Contig	ORF	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
18	6	5120	4218	gn P1D1d101118	YggG [Bacillus subtilis]	71	51	903
29	1	1	540	gi 1773142	similar to the 20.2kd protein in TEFB-200A region of B. subtilis [Bacillus subtilis]	71	56	540
38	20	113327	113830	gi 1537036	ORF_0158 [Escherichia coli]	71	48	504
51	12	15015	12676	gi 1109226	diacylglycerol lipase IV [Lactococcus lactis]	71	55	2140
55	23	124040	20585	gi 2143285	[AP015453] surface located protein [Lactobacillus rhamnosus]	71	58	456
60	2	705	265	gn P1D1d101320	YggG [Bacillus subtilis]	71	44	441
71	16	24679	26226	gi 560920	rodB (staA) polypeptide (AA 1-673) [Bacillus subtilis]	71	44	1548
71	25	30587	30360	gi 605028	ORF_0144; Geneplot suggests frameshift near start but none found [Bacillus subtilis]	71	50	228
72	6	5219	6729	gi 560835	lysine decarboxylase [Bacillus subtilis]	71	48	1491
72	14	11391	12878	gi 624085	414511; 50.5kDa heat-labile protease encoded by Cerebral Accession Number 827881; contains ATP/GTP binding motif [Paramecium bursaria Chlorella virus 1]	71	54	888
73	11	7249	7033	gi 1526559	Proteinectin (Protein)	71	42	237
74	6	10385	8517	gi 1573733	prolyl-tRNA synthetase (proS) [Haemophilus influenzae]	71	52	1869
81	9	4712	4578	gi 147406	Proteinase subunit II-4-Min [Escherichia coli]	71	45	807
86	5	4602	1604	gn P1D1e122063	ss-1,4-galactosyltransferase [Streptococcus pneumoniae]	71	53	599
105	4	3815	4707	gi 2223341	[AF014460] PspG [Streptococcus mutans]	71	58	1089
106	13	13557	12955	gi 1519287	LenA [Listeria monocytogenes]	71	48	603
114	2	1029	1979	gi 310361	moaA [Rhizobium meliloti]	71	55	951
122	2	564	1205	gi 1649017	glutamine transport ATP-binding protein GUNO [Salmonella typhimurium]	71	50	642
132	5	9018	7063	gn P1D1d102049	subtilisin	71	51	1356
140	1	1141	227	gi 1673788	[AC000015] Mycoplasma pneumoniae, fructose-bisphosphate aldolase; similar to Swiss-Prot Accession Number P13243, from B. subtilis [Mycoplasma pneumoniae]	71	49	915
140	5	5635	4973	gn P1D1d100964	homologue of hypothetical protein in a rapamycin synthesis gene cluster of Streptomyces hygroscopicus [Bacillus subtilis]	71	48	663
141	7	7389	7845	gn P1D1d102005	[AB001488] FUNCTION UNKNOWN; SIMILAR PRODUCT IN E. COLI AND MYCOPLASMA PNEUMONIAE. [Bacillus subtilis]	71	51	477

**TABLE 2**  
*S. pneumoniae* - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
193	1	1	165	[gi146912]	[ribosomal protein L13 (Staphylococcus carnosus)]	71	59	165
194	3	2205	1534	[gi153151]	[CofY (Bacillus subtilis)]	71	52	612
199	3	1310	1319	[gi12182574]	[A00009001 Y4pe (Rhizobium sp. NGR234)]	71	45	132
208	2	2616	3752	[gi11787718]	[A00002133 hypothetical protein in purB 5' region (Escherichia coli)]	71	57	1137
209	2	2022	1111	[gi141432]	[fepC gene product (Escherichia coli)]	71	46	882
210	5	1911	1071	[gi149316]	[ORF3 gene product (Bacillus subtilis)]	71	45	1161
210	6	3069	3386	[gi1580900]	[ORF3 gene product (Bacillus subtilis)]	71	48	318
212	2	3561	138	[gi155567]	[flonucleotide reductase H1 subunit (Mycobacterium tuberculosis)]	71	53	2181
233	3	2003	2920	[gm1ptd101020]	[yggB (Bacillus subtilis)]	71	50	918
244	1	13	1053	[gm1ptd1010944]	[homologue of aspartokinase 2 alpha and beta subunits LysC of B. subtilis (Bacillus subtilis)]	71	55	1041
251	2	1008	1874	[gi1755601]	[unknown (Bacillus subtilis)]	71	46	867
282	2	906	712	[gi11353874]	[unknown (Rhodococcus capsulatus)]	71	46	135
312	4	2137	1565	[gm1ptd1010245]	[A0005544 ynfB (Bacillus subtilis)]	71	34	573
338	1	3	683	[gi11571045]	[hypothetical protein (SP-P31466) (Methanococcus jannaschii)]	71	48	681
346	1	3	164	[gi11512124]	[hypothetical protein (SP-P4237) (Methanococcus jannaschii)]	71	36	162
374	1	619	2	[gi1379526]	[clumping factor (Staphylococcus aureus)]	71	23	618
377	1	688	2	[gi1379526]	[clumping factor (Staphylococcus aureus)]	71	23	687
3	8	7419	6958	[gm1ptd14263486]	[unknown (Bacillus subtilis)]	70	42	462
3	10	8395	9075	[gm1ptd14255543]	[putative iron dependent repressor (Staphylococcus epidermidis)]	70	46	681
7	14	11024	110254	[gm1ptd1010020]	[undefined open reading frame (Bacillus stearothermophilus)]	70	55	771
7	18	14213	13719	[gm1ptd1010090]	[biotin carboxyl carrier protein of acetyl-CoA carboxylase (Synchytrium sp.)]	70	56	495
9	2	1057	287	[gm1ptd1010081]	[unknown (Bacillus subtilis)]	70	52	771
12	4	2610	1789	[gm1ptd1010195]	[yycC (Bacillus subtilis)]	70	52	822
21	2	2586	1846	[gm1ptd1010195]	[ATPase (Bacillus subtilis)]	70	54	741
22	13	10955	11512	[gi11162295]	[Ydr340cp (Saccharomyces cerevisiae)]	70	50	558
30	6	4315	1980	[gi119478]	[ATP binding protein of transport ATPases (Bacillus firmus)]	70	51	336

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
31	1	370	113	gi1662792	single stranded DNA binding protein (unidentified bacterium)	70	36	258
33	15	10619	9521	gi11161219	homologous to D-amino acid dehydrogenase enzyme [Pseudomonas aeruginosa]	70	50	1119
38	6	3812	4312	gi1205847	ConrF [Streptococcus gordoni]	70	48	501
38	25	18477	18477	gi1517033	ORF_1556 [Escherichia coli]	70	58	492
40	113	11054	9846	gi11713516	riboflavin-specific deaminase [Actinobacillus pleuropneumoniae]	70	52	1209
42	2	722	1954	gi11146183	putative Bacillus subtilis	70	51	1233
43	3	3373	1612	gi11591493	glutamine transport ATP-binding protein 0 [Methanococcus jannaschii]	70	40	762
45	8	1897	8049	gn11701d100302	submit of ADP-glucose pyrophosphorylase [Bacillus stearothermophilus]	70	54	1149
59	2	567	956	gn11701d100302	neopullulanase [Bacillus sp.]	70	42	390
60	3	1874	795	gn11701d100302	aminopeptidase P [Lactococcus lactis]	70	48	1080
61	4	5553	2437	gn11701d100302	SNF [Bacillus cereus]	70	51	3177
61	7	7914	6802	gn11701d100302	cystathionine gamma-synthase [metB] [Haemophilus influenzae]	70	52	1113
63	7	5372	7232	gn11701d100974	unknown [Bacillus subtilis]	70	54	1851
68	7	7126	6942	gi1263014	emrA.1 gene product [Streptococcus pyogenes]	70	37	165
72	12	10081	10911	gi12313093	galactose-1-P-uridylyl transferase [Mycobacterium pyrolyticus]	70	56	831
75	10	7888	8124	gi11877423	galactose-1-P-uridylyl transferase [Streptococcus mutans]	70	59	237
79	3	3424	2525	gi139483	ORF 311 (AA 1-311) [Bacillus subtilis]	70	47	900
87	10	9369	7324	gn11701d100306	putative Pn2 protein [Bacillus subtilis]	70	52	2046
96	14	10440	11788	gi11573209	RNA-guanine transglycosylase (Gnt) [Haemophilus influenzae]	70	52	1149
113	2	574	1086	gi1433630	Al80 [Saccaromyces cerevisiae]	70	59	513
123	5	2901	3461	gn11701d100985	unknown [Bacillus subtilis]	70	45	561
125	5	4593	4282	gn11701d100985	capacitative calcium entry channel 1 [Bos taurus]	70	35	312
129	5	4500	3454	gn11701d101314	ORF [Bacillus subtilis]	70	47	1047
133	3	2608	1394	gi12233112	[AF008220] YcfP [Bacillus subtilis]	70	50	1215
135	1	420	662	gn11701d100530	YorF [Streptococcus pneumoniae]	70	47	243
137	3	438	932	gi1472919	v-type Na-ATPase [Enterococcus hirae]	70	57	495
138	3	440	3	gi1473316	transmembrane protein [Escherichia coli]	70	42	438

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
140	16	187896	16364	gi 976641	HS-methyltetrahydrofolate homocysteine methyltransferase (Saccharomyces cerevisiae)	70	53	2433
147	10	8263	6695	gi 149535	activating enzyme (Bacillus casei)	70	52	1569
204	4	3226	2747	gn FJ01d02049	E. coli hypochelator protein, p31805 (267) (Bacillus subtilis)	70	51	680
207	3	2627	2869	gn FJ01d09215	racGP (Dictyostelium discoideum)	70	45	243
282	3	1136	882	gi 1151874	unknown (Rhodospirillum rubrum)	70	50	235
6	21	11554	18453	gn FJ01d033079	hypochelator protein (Bacillus subtilis)	69	44	900
6	22	18482	18471	gi 560883	lpa-88d gene product (Bacillus subtilis)	69	53	990
22	6	4682	5834	gi 2205979	(AF006720) ProC (Bacillus subtilis)	69	48	1143
22	9	4992	8651	gn FJ01d00586	unknown (Bacillus subtilis)	69	51	660
22	12	3971	10797	gn FJ01d00581	unknown (Bacillus subtilis)	69	51	897
27	7	5857	5348	gn FJ01d02012	(AB001488) FUNCTION UNKNOWN (Bacillus subtilis)	69	53	2823
36	10	7284	10116	gi 437916	isoleucyl-tRNA synthetase (Staphylococcus aureus)	69	48	1089
38	1	2	1090	gi 141900	alcohol dehydrogenase (EC 1.1.1.1) (Alcaligenes eutrophus)	69	44	612
40	15	11333	11944	gi 1573280	Holliday junction DNA helicase (ruw) (Haemophilus influenzae)	69	50	576
40	15	11942	12517	gi 1573653	DNA-3-methyladenine glycosylase I (tagI) (Haemophilus influenzae)	69	47	1458
45	6	6947	5490	gi 580887	starch (bacterial glycogen) synthase (Bacillus subtilis)	69	36	780
48	34	24932	24153	gn FJ01d023870	hypochelator protein (Bacillus subtilis)	69	50	339
49	6	6183	6521	gi 3962297	similar to phosphotransferase system enzyme II (Escherichia coli)	69	49	733
49	8	7486	8338	gi 3964220	similar to Alcaligenes eutrophus pmi, D-ribulose-5-phosphate 3 epimerase (Escherichia coli)	69	50	1230
55	6	8262	7033	gn 1148238	poly(N) polymerase (Bacillus subtilis)	69	54	1380
59	3	934	2333	gn FJ01d01308	hypochelator protein (Bacillus subtilis)	69	49	249
62	3	1170	1418	gn FJ01d01915	hypochelator protein (Synechocystis sp.)	69	42	465
63	8	7298	7762	gi 293017	ORF3 (put.); putative (Lactococcus lactis)	69	49	1425
66	4	3657	5081	gi 153755	[phospho-beta-D-galactosidase (EC 3.2.1.45) (Lactococcus lactis cremoris)]	69	46	1704
66	5	5126	4829	gi 433809	enzyme II (Streptococcus mutans)	69	39	648
71	6	10017	10664	gn FJ01d022063	ss-1,4-galactosyltransferase (Streptococcus pneumoniae)	69	49	733

TABLE 2

S. pneumoniae - Relative coding regions of novel proteins similar to known proteins

Contig ID	Start ID	Stop ID	match accession	match gene name	% sim	% ident	length (nt)
71	21	27730	gnt1P1D100649	DE-cadherin [Drosophila melanogaster]	69	30	237
77	1	237	gnt1P1D100649	gnt1 gene product [Lactococcus lactis]	69	44	237
81	5	3622	gnt1P1D100649	fucose operon protein (fucP) [Haemophilus influenzae]	69	52	480
83	16	15742	gnt1P1D100649	phosphoribosyl glycineamide formyltransferase (ribA) [Bacillus subtilis]	69	46	675
85	2	1212	gnt1P1D100649	IFN-response element binding factor 1 [Mus musculus]	69	46	594
91	5	3678	gnt1P1D100649	anaerobic ribonucleoside-triphosphate reductase activating protein (nrdr) [Haemophilus influenzae]	69	44	597
98	5	3247	gnt1P1D100649	lyf protein [Salmonella typhimurium]	69	51	786
108	5	4085	gnt1P1D100649	transcription factor [Lactococcus lactis]	69	49	972
126	3	1078	gnt1P1D100649	Yqj [Bacillus subtilis]	69	49	1491
131	6	4121	gnt1P1D100649	Yqk [Bacillus subtilis]	69	47	1233
136	2	2525	gnt1P1D100649	unknown [Bacillus subtilis]	69	47	795
149	5	3852	gnt1P1D100649	VioQ protein [Bacillus subtilis]	69	50	912
149	12	9336	gnt1P1D100649	homology with E.coli and P.aeruginosa lytA gene; product of unknown function [Pseudomonas syringae]	69	52	1320
153	4	3191	gnt1P1D100649	BrnQ [Bacillus subtilis]	69	44	639
169	3	849	gnt1P1D100649	temperature sensitive cell division [Bacillus subtilis]	69	49	1476
180	1	566	gnt1P1D100649	alpha-amyase [unidentified cloning vector]	69	50	544
212	1	1196	gnt1P1D100649	ribonucleoside reductase R2-2 small subunit [Mycobacterium tuberculosis]	69	53	966
226	1	2	gnt1P1D100649	modulin-26 - soybean	69	41	660
231	5	2240	gnt1P1D100649	gamma-type Na-ATPase [Enterococcus hirae]	69	56	1518
235	3	660	gnt1P1D100649	methyase [Haemophilus influenzae]	69	43	1107
253	2	865	gnt1P1D100649	ORF5 [Barley yellow dwarf virus]	69	69	1497
251	3	2899	gnt1P1D100649	neurolefin-efflux protein [Stenotrophomonas maltophilia]	69	51	933
310	1	282	gnt1P1D100649	peptide deformylase [Clostridium botulinum]	69	55	242
369	1	868	gnt1P1D100649	clumping factor [Staphylococcus aureus]	69	22	867
370	1	749	gnt1P1D100649	clumping factor [Staphylococcus aureus]	69	21	747

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start	Stop	Match accession	Match gene name	% sim	% ident	Length (nt)
379	1	44	280	gi 171010649	DK-ectherin (Drosophila melanogaster)	69	30	237
388	1	280	72	gi 1787524	[AB000225] hypothetical 32.7 kb protein in trpA-btuR intergenic region (Escherichia coli)	69	44	189
1	2	2066	3040	gi 17101809	ABC transporter (Synchocystis sp.)	68	43	1035
12	5	3958	2600	gi 2182492	histidine kinase [Lactococcus lactis cremoris]	68	45	1359
15	2	1790	1311	gi 1619741888	ribosomal protein L9 - Bacillus steatothermophilus	68	50	460
16	6	7353	5701	gi 1787041	[AB000144] c530, this 530 aa orf is 33 pct identical (41 gaps) to 525 residues of an approx. 840 aa protein YMS_J0429 SM: F44808 [Escherichia coli]	68	45	1653
17	12	1479	4805	gi 1553165	acetylcholinesterase (Homo sapiens)	68	68	327
20	13	14128	14505	gi 142700	P competence protein (ttg start codon) (put.), putative [Bacillus subtilis]	68	40	378
22	12	24632	25377	gi 126262	comE ORF3 [Bacillus subtilis]	68	36	786
30	7	4548	4288	gi 311388	ORF1 [Asorhizobium caulinodans]	68	46	261
36	5	3211	4585	gi 1571041	hypothetical [Memphilius influenzae]	68	54	675
46	6	5219	6040	gi 1790131	[AB000446] hypothetical 39.7 kD protein in lpaA-yyrB intergenic region [Escherichia coli]	68	47	822
54	10	6235	7086	gi 1882579	CO 6180 aa c3759 [Escherichia coli]	68	55	852
55	5	7069	5165	gi 17101914	ABC transporter [Synchocystis sp.]	68	45	1905
71	3	6134	5633	gi 1573355	outer membrane integrity protein (tolA) [Memphilius influenzae]	68	50	522
71	10	15342	16613	gi 1580866	lpa-12d gene product [Bacillus subtilis]	68	31	1272
71	12	17360	18792	gi 144075	SecY protein [Lactococcus lactis]	68	35	1233
71	17	22295	24703	gi 1762149	Involved in protein export [Bacillus subtilis]	68	50	2409
73	16	10208	9729	gi 1353537	DUF696 [Bacteriophage phi]	68	51	480
86	18	17198	18011	gi 1413943	lpa-13d gene product [Bacillus subtilis]	68	53	1188
87	17	17494	15866	gi 150209	ORF 1 [Mycothecium mycoides]	68	43	1656
89	6	5139	4354	gi 1458824	M. jamaicensis predicted coding region K4062 [Lactococcus jamaicensis]	68	40	786
89	11	8021	8242	gi 150974	[4-oxalacetate] tautomerase [Pseudomonas putida]	68	43	222
97	8	6755	5394	gi 2367358	[AB000491] hypothetical 52.9 kD protein in aldA-rpsI intergenic region [Escherichia coli]	68	41	1362

TABLE 2  
S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match gene name	% sim	% ident	Length (nt)
98	3	1418	2308	gi 1100261	[Lys protein (Salmonella typhimurium)	68	40	891
99	13	16414	17280	gi 455363	transitory protein (Streptococcus mitis)	68	50	867
115	3	3054	3693	gi 446474	cellulose phosphotransferase enzyme I <sup>+</sup> (Bacillus thermoautotrophicus)	68	44	1362
124	7	3394	321	gi 10100702	core protein (Schistosoma mansoni)	68	56	171
125	2	2923	1922	gi 420566	transmembrane protein (Bacillus subtilis)	68	50	1002
132	2	4858	2888	gi 10100702	DNA ligase (Synchococcus sp.)	68	52	1271
140	7	7745	7580	gi 1209711	unknown (Saccharomyces cerevisiae)	68	47	186
150	1	539	3	gi 402490	put-ribosylarginine hydrolase (Mus musculus)	68	59	537
164	1	58	867	gi 10100702	glutamate racemase (Bacillus subtilis)	68	49	810
164	2	819	1835	gi 10100702	hypothetical protein (Bacillus subtilis)	68	50	1017
169	7	3946	4104	gi 10100702	hypothetical protein - Lactococcus lactis subsp. lactis plantarum	68	40	159
170	4	4247	4396	gi 10100702	spore coat protein (Bacillus subtilis)	68	52	150
171	8	6002	7054	gi 10100702	precursor (aa -20 to 181) (Bacillus subtilis)	68	54	1053
198	3	2473	1871	gi 10100702	hypothetical protein (Bacillus subtilis)	68	46	603
211	2	969	1802	gi 1439828	BTIC-man (Bacillus subtilis)	68	45	834
214	8	4926	4231	gi 10100702	H. influenzae hypothetical protein; P43990 (B2) (Bacillus subtilis)	68	50	696
217	6	4955	5170	gi 10100702	similar to 8 volutinase associated mitochondrial ... (reverse transcriptase) (Haeberlein coli)	68	36	216
218	7	3930	4745	gi 1225199	(AF008220) Ycp (Bacillus subtilis)	68	38	816
220	6	4628	4338	gi 10100702	(AF008005) orf1 (Bacillus megaterium)	68	51	291
236	1	746	108	gi 1010177	ORF13 (Bacillus subtilis)	68	46	639
237	2	675	1451	gi 1396348	homoserine transaminase (Bacillus coli)	68	49	737
250	4	771	1229	gi 1010859	ORF2 (Synchococcus sp.)	68	50	459
254	1	517	155	gi 1787105	(AF000189) orf4 was 6649, this 669 aa orf is 40 pct identical (1 gap) to 217 residues of an approx. 232 aa protein YBA <sub>10</sub> (GenBank: F45247) (Bacillus coli)	68	44	363
337	1	1	774	gi 10100702	putative orf (Bacillus subtilis)	68	47	774
345	3	653	1455	gi 1495513	thymidylate synthase (EC 2.1.1.4) (Lactococcus lactis)	68	61	651



TABLE 2  
S. pneumoniae - Native coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% Ident. (nt)	length (nt)
386	2	417	4	gi11571353	outer membrane integrity protein (toxA) [Haemophilus influenzae]	68	51	114
2	4	5722	4697	gi11592141	M. jannaschii predicted coding region M1507 [Methanococcus jannaschii]	67	26	1028
3	6	5397	4591	gi12293175	(AF008220) signal transduction regulator [Bacillus subtilis]	67	44	107
5	2	2301	574	gi12131385	(AF000547) para-aminobenzoate synthetase (pabI) [Helicobacter pylori]	67	48	1728
6	19	14643	16758	gi1413931	Ipa-td gene product [Bacillus subtilis]	67		
22	8	7024	7897	gi12928962	pyruvate 5-carboxylate reductase [Methanobrevibacter smithii]	67	51	804
29	10	8335	9072	gi1468745	ptcB gene product [Bacillus brevis]	67	41	728
31	2	1379	585	gi12425123	(AF019986) PsaB [Dictyostelium discoideum]	67	49	795
32	11	8849	10150	gi1420029	ORF1 gene product [Escherichia coli]	67	47	1102
36	16	14830	15546	gi11921442	ABC transporter, probable ATP-binding subunit [Methanococcus jannaschii]	67	43	717
38	9	4938	4392	gi1170141803	P2333.3 [Methanobrevibacter smithii]	67	47	435
38	121	11775	14512	gi1527037	ORF_0216 [Escherichia coli]	67	52	738
45	9	10428	9181	gi1541710	...-chain enzyme (gipB) [EC 2.4.1.18] [Bacillus acetotolerophilus]	67	51	1248
48	123	18344	17514	gi1413949	Ipa-5d gene product [Bacillus subtilis]	67	50	831
50	2	1773	952	gi117014101320	Yqj00 [Bacillus subtilis]	67	55	822
53	1	431	3	gi11574293	fibribral transcription regulation repressor (pilB) [Haemophilus influenzae]	67	40	425
55	13	12740	11946	gi112101425990	ORF T00037c [Saccharomyces cerevisiae]	67	51	795
61	9	9210	8329	gi1171014264711	ATP-binding cassette transporter A [Staphylococcus aureus]	67	50	882
71	2	3814	617	gi11197867	vitellogenin [Anolis pulchellus]	67	36	504
81	7	1489	4583	gi11142714	phosphoenolpyruvate-mannose phosphotransferase element 118 [Lactobacillus curvatus]	67	42	495
83	7	2857	3214	gi12767646	...-carrier protein [Porphyra purpurea]	67	37	258
86	8	8140	6809	gi11147744	PER [Enterococcus hirae]	67	45	1332
97	3	386	356	gi117101402235	(AB006031) unanased protein product [Streptococcus mutans]	67	43	381
102	1	601	1413	gi1682785	mecB gene product [Escherichia coli]	67	36	813
106	3	2607	1907	gi1146921	licP protein [Haemophilus influenzae]	67	43	879
115	4	5982	5656	gi1895750	putative cellulose phosphotransferase enzyme III [Bacillus subtilis]	67	44	327

TABLE 2  
5. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	OHF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
115	7	8421	8077	gi1464473	cellulose phosphotransferase enzyme II <sup>+</sup> [Bacillus stercorophilus]	67	51	345
127	13	8127	7021	gi147326	transport protein [Escherichia coli]	67	45	1107
136	3	2215	2859	gn1P1D1d100581	unknown [Bacillus subtilis]	67	49	645
140	21	23317	20506	gn1P1D1d101912	phenylalanyl-tRNA synthetase [Synchocystis sp.]	67	43	2412
146	6	2884	1893	gn12182994	histidine kinase [Lactococcus lactis cremoris]	67	44	1002
151	8	11476	11117	gn1P1D1d100085	[ORF129 [Bacillus cereus]	67	48	360
160	10	7451	6646	gi12281317	[ORF6], similar to a Streptococcus pneumoniae putative membrane protein encoded by GenBank Accession Number X591400; inactivation of the Orf6 gene leads to UV-sensitivity and to decrease of homologous recombination (plasmidic test) [Lactococcus 1]	67	46	1194
163	3	3099	4505	gn1P1D1d101317	YnfR [Bacillus subtilis]	67	47	1407
167	8	6704	5454	gi1161933	[DltB [Lactobacillus casei]	67	45	1251
169	4	2322	2879	gn1P1D1d101331	YnfG [Bacillus subtilis]	67	41	556
171	11	7656	8384	gi1513841	pneumococcal surface protein A [Streptococcus pneumoniae]	67	50	729
188	3	1930	1733	gi1542735	DltB [Mycobacterium thermophilus]	67	46	1794
189	6	3599	3141	gn1P1D1e35178	hypothetical protein [Bacillus subtilis]	67	52	459
205	3	366	2211	gi1060675	[ORF_0169 [Escherichia coli]	67	47	549
207	4	2896	3456	gi12276374	DltB/Iron regulated lipoprotein precursor [Corynebacterium diphtheriae]	67	49	561
217	3	4086	3703	gi1095750	putative cellobiose phosphotransferase enzyme III [Bacillus subtilis]	67	42	384
246	2	291	662	gi11842438	unknown [Bacillus subtilis]	67	43	372
252	1	2	745	gi12351766	PapA [Streptococcus pneumoniae]	67	41	744
265	3	1134	1811	gn12313487	(A00005185) L-asparaginase II (anab) [Helicobacter pylori]	67	42	678
295	1	1	375	gi12276374	DltB/Iron regulated lipoprotein precursor [Corynebacterium diphtheriae]	67	43	375
3	7	4898	5146	gn1P1D1e255179	unknown [Mycobacterium tuberculosis]	66	56	249
3	1	389	3	gn1P1D1e269548	unknown [Bacillus subtilis]	66	48	387
3	20	19267	20805	gi139956	[DltC [Bacillus subtilis]	66	50	1539
4	3	2545	2718	gi11787564	(A0000218) phase shock protein C [Escherichia coli]	66	36	174
5	9	13197	12592	gi1574291	[GmbA1 transcription regulation repressor (p118) [Haemophilus influenzae]	66	46	606

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession (nt)	match gene name	% sim	% ident	length (nt)
9	4	2872	1451	gi 121026628	unknown [Mycobacterium tuberculosis]	66	43	1122
12	2	1469	1200	gi 1520407	orf2: ORF start codon [Bacillus thuringiensis]	66	42	270
15	12	10979	9897	gi 2314398	(AK006653) translation elongation factor EF-Ts (tsf) [Helicobacter pylori]	66	49	1093
16	2	1312	734	gi 1402245	(AK005554) y2b [Bacillus subtilis]	66	35	579
22	3	1372	1851	gi 1460916	signal peptidase type II [Lactococcus lactis]	66	38	480
22	7	5828	7096	gi 121026628	gamma-glutamyl phosphate reductase [Streptococcus thermophilus]	66	51	1269
22	20	16194	17138	gi 1210268191	Y12L [Bacillus subtilis]	66	50	945
30	2	530	976	gi 2314379	(AK006677) ABC transporter, ATP-binding protein (yhqG) [Helicobacter pylori]	66	40	447
32	1	199	984	gi 312444	ORF2 [Bacillus caldolyticus]	66	49	786
33	13	8352	7234	gi 1387979	44% identity over 302 residues with hypothetical protein from Synchocystis sp. PCC 6803. The protein is involved by iron homeostasis, some similarity to glycyl transferases, two potential membrane-spanning helices [Bacillus subtilis]	66	44	1119
34	6	5458	1708	gi 121026724	orf2 [Mycobacterium avium]	66	39	951
34	14	9792	9574	gi 1590597	M. jannaschii predicted coding region M20772 [Methanococcus jannaschii]	66	48	219
35	16	15180	14503	gi 1773532	[CapSM (Staphylococcus aureus)]	66	46	663
36	9	6173	6976	gi 1518680	Minical-associated protein DivIVA [Bacillus subtilis]	66	35	894
36	11	10396	10874	bbi 155344	Insulin activator factor, INSAF (human, Pancreatic insulinoma, Peptide hormone) [Homo sapiens]	66	43	479
48	1	28	1419	gi 1210262504	hypothetical protein [Bacillus subtilis]	66	50	1392
48	7	3810	4112	gi 2182874	(AK006090) Y4b [Bacillus sp. 5023141]	66	40	303
52	4	3395	2789	gi 1388565	major cell-binding factor [Campylobacter jejuni]	66	52	867
54	3	2662	1076	gi 1210261031	cell-binding periplasmic protein [Synchocystis sp.]	66	43	1587
61	10	9740	9183	gi 1210261444	hdr gene product [Staphylococcus aureus]	66	44	558
72	3	10895	11993	gi 2313129	(AK006026) H. pylori predicted coding region HP0049 [Helicobacter pylori]	66	44	1101
74	9	13267	12476	gi 1573941	hypothetical [Methanophilus influenzae]	66	43	792
75	1	2	868	gi 1574631	nicotinamide mononucleotide transporter (pncC) [Methanophilus influenzae]	66	48	867
75	7	5103	4275	gi 41312	put. BRG repressor protein [Escherichia coli]	66	40	1029



TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
223	1	1070	118	[gm PRDj247187	zinc finger protein [Bacteriophage phigle]	66	45	933
224	1	1864	2640	[gm 1176199	putative ABC transporter subunit [Staphylococcus epidermidis]	66	41	777
243	1	3	972	[dbj JAO00617.2	[AB000617] YcdH [Bacillus subtilis]	66	45	870
268	2	891	168	[gi 517210	putative transposase [Streptococcus pyogenes]	66	40	324
32	1	2	643	[gi 1499836	[zn protease [Methanococcus jannaschii]	66	40	642
5	10	13909	11178	[gi 1574292	hypothetical [Hemophilus influenzae]	65	34	732
6	11	10465	11190	[gi 142854	homologous to E. coli rdc gene product and to unidentified protein from [Staphylococcus aureus [Bacillus subtilis]	65	48	726
7	2	647	405	[pir C64146 C641	hypothetical protein NI0259 - Hemophilus influenzae (strain Rd RM20)	65	42	243
7	3	5246	6821	[gm PRDj010123	[gdu [Bacillus subtilis]	65	50	576
10	2	1873	1397	[gi 1163111	[ORF-1 [Streptococcus pneumoniae]	65	54	477
16	3	1426	2222	[gm PRDj25010	hypothetical protein [Bacillus subtilis]	65	45	795
21	4	3815	3357	[gm PRDj314910	hypothetical protein [Staphylococcus sciuri]	65	40	459
22	34	25776	26384	[gi 1123030	[opa [Actinobacillus pleuropneumoniae]	65	42	609
43	2	1648	280	[gi 1044826	[F1485.1 [Ctenocephalides felis]	65	38	1359
46	13	10962	10486	[gi 1573391	hypothetical [Hemophilus influenzae]	65	45	795
48	22	17321	16883	[gi 1573391	hypothetical [Hemophilus influenzae]	65	37	639
46	25	19027	18553	[gm PRDj254404	[CR020C. Ien:215 [Saccharomyces cerevisiae]	65	30	495
49	3	3356	5334	[gi 1480428	putative transcriptional regulator [Bacillus stearothermophilus]	65	32	1179
50	6	5337	4519	[gi 171963	[RNA isopentenyl transferase [Saccharomyces cerevisiae]	65	42	819
52	15	14728	15588	[gi 1499745	[M. jannaschii predicted coding region M3912 [Methanococcus jannaschii]	65	46	861
59	7	3363	4745	[gi 456514	[orf zeta [Streptococcus pyogenes]	65	42	783
68	3	2501	3483	[gi 987824	[omp a310 [Escherichia coli]	65	46	984
69	3	2171	1077	[gm PRDj311453	unknown [Bacillus subtilis]	65	42	1095
69	7	6029	5325	[gi 1095660	unknown [Bacillus subtilis]	65	55	705
71	5	8536	9783	[gi 1573224	[glycyl transferase IgtC (GP-016584.4) [Hemophilus influenzae]	65	42	1288
72	8	7664	8527	[gm PRDj257049	unknown, highly similar to several spermidine synthases [Bacillus subtilis]	65	39	864

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	length (aa)
76	5	5715	6977	gnl pjd g101723	DNA REPAIR PROTEIN REC1 (RECOMBINATION PROTEIN N.) [Escherichia coli]	65	44	1677
76	9	8099	7875	gi 1574276	acetoacetylcholinesterase, small subunit (scab) [Haemophilus influenzae]	65	38	225
84	2	2870	2152	gi 1231188	[AB000532] conserved hypothetical protein [Helicobacter pylori]	65	41	519
86	15	14495	13407	gnl pjd g101880	3-dehydroquinate synthase [Synecococcus sp.]	65	44	1089
87	3	2465	2423	gi 151259	HMG-CoA reductase (EC 1.1.1.88) [Pseudomonas mewanii]	65	51	1284
88	3	3706	2736	gi 1098510	unknown [Lactococcus lactis]	65	30	312
89	2	1627	1007	gnl pjd g102008	[AB001488] SIMILAR TO ORF14 OF ENTEROCOCCUS FAECALIS TRANSFOSIN T916. [Bacillus subtilis]	65	41	621
111	6	6635	6186	gnl pjd g246063	HK23/nucleoside diphosphate kinase [Monopus laevis]	65	50	450
116	1	730	1016	gnl pjd g101125	Queuosine biosynthesis protein QueA [Synecococcus sp.]	65	44	1014
123	1	69	189	gi 198839	ORF2 [Clostridium perfringens]	65	36	321
123	7	6322	7190	gi 1575577	DNA-binding response regulator [Thermotoga maritima]	65	39	669
125	3	3921	2859	gnl pjd g257609	sugar-binding transport protein [Anaerococcus thermophilus]	65	47	965
137	12	8015	7818	gi 1262574	[AB000090] Y4p [Rhizobium sp. NOR234]	65	41	198
147	4	5021	3885	gi 172329	dihydrolipoamide acetyltransferase [Clostridium magnum]	65	47	1137
148	2	1055	1931	gnl pjd g101319	YggH [Bacillus subtilis]	65	42	879
151	2	3212	4687	gi 1304897	Ecse type 1 restriction modification enzyme N subunit [Escherichia coli]	65	50	1476
156	2	730	437	gi 1310893	membrane protein [Thalassia parva]	65	47	294
164	7	4256	4837	gi 1410132	ORF28 [Bacillus subtilis]	65	48	582
169	6	3192	3914	gi 1552777	similar to purine nucleoside phosphorylase (deob) [Escherichia coli]	65	41	723
176	4	2951	2220	gnl pjd g319500	oligonucleotide binding lipoprotein [Streptococcus pneumoniae]	65	43	732
195	4	4556	3900	gi 1592142	ABC transporter, probable ATP-binding subunit [Methanococcus jannaschii]	65	40	657
196	1	160	1572	gnl pjd g102004	[AB001489] PROBABLY NON-ATP-BINDING SUBUNIT OF A DUMP/ATL-2, 6- OXIMINOLIPASE (EC 6.3.1.5) [Bacillus subtilis]	65	51	1413
204	2	2246	1215	gi 1431356	membrane bound protein [Bacillus subtilis]	65	37	1032
210	4	1544	1851	gi 149315	ORF7 gene product [Bacillus subtilis]	65	48	348
242	2	1625	723	gi 1787540	[AE000226] f249; This 249 aa orf is 32 pct identical (8 gaps) to 244 residues of an approx. 272 aa protein ANA-BCOI SW: P42903 [Escherichia coli]	65	42	903

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
284	1	1	900	[gi 559861	clpX (Plasmid pAD1)	65	36	900
304	1	2	574	[gn P1D 029034	unknown [Mycobacterium tuberculosis]	65	52	573
315	1	2	1483	[gi 790654	hemucron C-5-epimerase [Acetobacter vinelandii]	65	57	1482
320	1	3	569	[gn P1D 02048	K. aerogenes, histidine utilization repressor, P12300 [199] DNA binding [Bacillus subtilis]	65	46	567
358	1	309	[gn P1D 023508	ProG protein [Bacillus subtilis]		65	55	309
2	7	7571	6656	[gi 1489753	nicotinate-nucleotide pyrophosphorylase [Rhodospirillum rubrum]	64	47	876
6	6	5924	4802	[gn P1D 010111	histonone aminopeptidase [Synecocystis sp.]	64	52	879
8	4	3417	3686	[gi 1045935	DNA helicase II [Mycoplasma genitalium]	64	58	270
11	4	3449	2689	[gn P1D 0265539	OrfB [Streptococcus pneumoniae]	64	46	561
15	7	6504	7245	[gi 1762328	Yer59c/YigZ homolog [Bacillus subtilis]	64	45	642
22	11	9348	9895	[gn P1D 0100581	unknown [Bacillus subtilis]	64	38	348
22	30	22593	23174	[gi 289260	[cong ORF1 [Bacillus subtilis]	64	44	672
26	7	14275	14199	[gi 405286	[berru [Bacillus subtilis]	64	30	177
27	2	1510	1334	[gi 40795	[cdaI methylase [Desulfovibrio vulgaris]	64	51	177
29	2	614	297	[gi 226168	[cype VII collagen [Mus musculus]	64	50	318
35	2	368	721	[pir JC1153 JC11	hypothetical 20.3K protein (insertion sequence IS1131) - Agrobacterium tumefaciens (strain P022) plasmid T1	64	50	334
40	1	3	449	[gi 46970	[spid gene product [Staphylococcus epidermidis]	64	41	447
40	7	4683	4976	[gn P1D 0325792	[AJ000005] glucose kinase [Bacillus megaterium]	64	45	234
45	7	8068	6920	[gn P1D 020036	subunit of ADP-glucose pyrophosphorylase [Bacillus stearothermophilus]	64	40	1149
51	2	301	1059	[gi 43985	[nifs-like gene [Acetobacillus delbrueckii]	64	54	759
51	13	15551	10397	[gi 2293260	[AF002220] DNA-polymerase III alpha-chain [Bacillus subtilis]	64	46	3147
53	3	1157	555	[gi 1574822	[hypothetical [Mycobacterium influenzae]	64	47	603
58	2	4236	1606	[gi 1573826	[alanyl-tRNA synthetase (ala) [Mycobacterium influenzae]	64	51	2631
66	1	3	1259	[gi 655749	[positive collagenase phosphotransferase enzyme 11'' [Bacillus subtilis]	64	42	1257
68	5	5213	6556	[gi 416965	[metA] gene products [Bacillus stearothermophilus]	64	47	1344
69	6	5356	4949	[gn P1D 010316	[cdd [Bacillus subtilis]	64	52	408

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Gene ID	ORF ID	Start (nt)	End (nt)	Accession	Match gene name	% ident	% ident	Length (nt)
74	4	6948	5038	gi 726480	l-glutamine-D-fructose-6-phosphate amidotransferase [Bacillus subtilis]	64	50	1911
75	3	1283	1465	bbai133179	YUS-OmpA fusion protein(OmpA/OMP transcription factor, TUS-nuclear RNA-binding protein) [human, myxoid liposarcoma cells, peptide mutant, 482 aa] [Homo sapiens]	64	57	183
81	13	14016	14231	gi 143175	methanol dehydrogenase alpha-10 subunit [Bacillus sp.]	64	35	216
83	22	21851	12090	gnl pfj d101315	Yqda [Bacillus subtilis]	64	44	240
87	11	10946	9100	gnl pfj d23505	putative PcaI protein [Bacillus subtilis]	64	43	747
98	7	5032	2706	gnl pfj d23880	hypothetical protein [Bacillus subtilis]	64	38	675
105	3	2	1276	gnl pfj d101314	similar to S. pneumoniae Yqda	64	45	1275
113	7	5136	4410	gnl pfj d101119	NifS [Synecocystis sp.]	64	50	1275
118	3	1287	1750	gnl pfj d20520	hypothetical protein [Macromonococcus parvulus]	64	37	1286
123	3	1125	2156	gnl pfj d23284	ORF Yqda244w [Saccharomyces cerevisiae]	64	40	1032
124	5	2331	1780	gnl pfj d101844	hypothetical protein [Synecocystis sp.]	64	50	552
129	4	3467	2709	gnl pfj d101314	Yqda [Bacillus subtilis]	64	52	759
131	3	122	3	gnl pfj d101314	unknown [Bacillus subtilis]	64	42	150
137	11	7196	7549	pir JC1151 3C11	hypothetical 20.3K protein (insertion sequence IS1131) - Agrobacterium tumefaciens (strain P022) plasmid T1	64	50	354
139	3	3225	2651	gi 2353303	(AF004220) Yqda [Bacillus subtilis]	64	44	576
146	10	6730	5648	gi 1322245	isovalerate pyrophosphate decarboxylase [Rattus norvegicus]	64	45	1083
147	3	2	1018	gnl pfj d237033	unknown gene product [Acetobacterium leifmannii]	64	46	1017
148	11	8430	8783	gi 2110630	(AF0004301) dynamin-like protein [Homo sapiens]	64	28	354
156	7	4313	3312	gnl pfj d100250	transmembrane [Bacillus subtilis]	64	31	702
157	4	1299	2114	gnl pfj d100892	homologous to Gln transport system permease proteins [Bacillus subtilis]	64	43	816
162	6	5880	6182	gnl pfj d101204	ORF1, putative 42 kDa protein [Streptococcus pyogenes]	64	58	483
164	13	9707	8769	gnl pfj d100964	homologue of ferric anquibactin transport system permease protein PATD of V. anguillarum [Bacillus subtilis]	64	40	939
175	5	3905	4598	gi 524045	actin [Bacillus subtilis]	64	39	693
189	10	6156	6507	gnl pfj d101307	response regulator [Acetobacterium plantarum]	64	33	354
191	4	2539	2465	gnl pfj d100520	phosphotransferase [Lactococcus lactis]	64	46	657



TABLE 2  
S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% Ident	length (nt)
202	1	76	1140	gi 1293806	[p-actylhomoserine sulphydriase (Streptococcus agalactiae)]	64	47	1065
224	1	234	1571	gi 1571393	[collagenase (prtC) (Haemophilus influenzae)]	64	42	1338
231	3	291	647	gi 140174	[ORF X (Bacillus subtilis)]	64	43	357
253	3	709	1089	gi 137151 JC11	[hypothetical 20.3K protein (insertion sequence IS1131) - Agrobacterium tumefaciens (strain P02) plasmid Ti]	64	50	381
265	1	820	2	gi 1177832	[unknown (Bacillus subtilis)]	64	31	812
297	1	1	660	gi 1590871	[collagenase (Methanococcus jannaschii)]	64	48	660
328	1	263	21	gi 192651	[GlnP (Saccharomyces cerevisiae)]	64	41	243
5	4	8730	8998	gi 1556685	[unknown (Bacillus subtilis)]	63	48	633
10	6	5178	4483	gi 1571101	[hypothetical (Haemophilus influenzae)]	63	40	696
12	11	9324	9502	gi 1806536	[membrane protein (Bacillus acidopullulicus)]	63	42	579
15	110	8897	9187	gi 172219	[unknown (Acetobacter xylinum)]	63	40	291
17	2	1031	309	gi P1D e217602	[p1nd (Lactobacillus plantarum)]	63	32	723
18	8	7778	6975	gi 1177843	[unknown (Bacillus subtilis)]	63	45	804
26	4	9780	7078	gi 142440	[ATP-dependent nuclease (Bacillus subtilis)]	63	46	2703
29	5	3488	4192	gi 1177829	[unknown (Bacillus subtilis)]	63	35	705
34	11	8830	7988	gi P1D d101198	[ORF8 (Enterococcus faecalis)]	63	45	843
35	2	1120	876	gi 1722339	[unknown (Acetobacter xylinum)]	63	39	312
48	15	12509	11691	gi 1571389	[hypothetical (Haemophilus influenzae)]	63	41	819
51	11	12719	5022	gi 142450	[ORF3 protein (Bacillus subtilis)]	63	35	531
55	4	3979	5022	gi 1708640	[YsaB (Bacillus subtilis)]	63	41	1044
55	15	13669	14670	gi P1D e31502	[thioredoxine reductase (Bacillus subtilis)]	63	44	1002
68	10	9242	8919	gi P17686 V1A7	[HYPOPHOSPHAL 40.2 KD PROTEIN IN AV19-SERIE INTERGENIC REGION (P1822)]	63	40	324
86	7	5254	5685	gi 1571382	[Ile-1 operon protein (IleD) (Haemophilus influenzae)]	63	41	870
88	8	6085	5180	gi 2098719	[putative fibriin-associated protein (Actinomyces viscosus)]	63	43	566
96	8	2858	6484	gi 1052803	[orf19yb gene product (Streptococcus pneumoniae)]	63	38	627
100	1	240	1940	gi 1711	[fucosidase (Dictyostelium discoideum)]	63	36	1701

TABLE 2

S. pneumoniae - putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
104	4	3063	5765	gi 144985	phosphoenolpyruvate carboxylase (Clostridium glutamicum)	63	46	2703
106	8	9189	8554	gi 533099	endonuclease III (Bacillus subtilis)	63	45	636
122	6	4704	4886	gi 190101139	transase [Synecocystis sp.]	63	39	183
128	7	4517	5203	gi 190101434	orf2 (Methanobacterium thermoautotrophicum)	63	50	687
137	4	963	1547	gi 1722920	Type III-A kinase (Enterococcus hirae)	63	27	585
142	7	4100	4385	gi 190131025	hypothetical protein [Bacillus subtilis]	63	44	486
159	5	1741	2571	gi 1187043	hypothetical protein [Z. m. 213] or 14-24 pct identical (18 gaps) to 265 residues of an approx. 272 aa protein Y100_0001 [Escherichia coli]	63	39	831
171	112	8803	14406	gi 1901324918	galactose 4-epimerase [Streptococcus sanguis]	63	48	5604
177	1	3	347	gi 1773150	hypothetical 14.8kd protein [Escherichia coli]	63	34	345
178	2	473	517	gi 1722335	unknown [acetobacter xylinum]	63	41	495
178	3	794	1012	gi 1591582	cobalamin biosynthesis protein N [Methanococcus jannaschii]	63	36	219
195	1	1377	175	gi 1901324217	ftsQ [Enterococcus hirae]	63	33	1203
234	5	1739	1527	gi 1591582	cobalamin biosynthesis protein N [Methanococcus jannaschii]	63	36	213
249	1	81	257	gi 1000453	TrkA [Bacillus subtilis]	63	41	177
283	1	127	1347	gi 1396486	ORF8 [Bacillus subtilis]	63	44	1221
293	3	2804	3466	gi 1722339	unknown [acetobacter xylinum]	63	37	663
311	1	905	486	gi 1877424	UDP-galactose 4-epimerase [Streptococcus mutans]	63	46	420
324	1	2	556	gi 1477741	histidine periplasmic binding protein P29 [Campylobacter jejuni]	63	36	555
365	1	219	13	gi 1255843	(AF013293) no definition line found [Methanobacterium jannaschii]	63	33	207
382	1	88	378	gi 1722339	unknown [acetobacter xylinum]	63	40	291
382	3	364	158	gi 1255843	(AF013293) no definition line found [Methanobacterium jannaschii]	63	33	207
2	1	2495	288	gi 1901329007	penicillin-binding protein [Bacillus subtilis]	62	42	2208
3	23	12374	12421	gi 1901325493	hypothetical protein [Bacillus subtilis]	62	35	858
6	16	14320	13193	gi 1901324914	in/f5-like protein [Mycobacterium leprae]	62	37	1128
7	8	6819	7322	gi 190101324	orf1 [Bacillus subtilis]	62	32	414
7	19	15166	14207	gi 190101804	beta ketacyl acyl carrier protein synthase [Synecocystis sp.]	62	43	1260

TABLE 2  
5. pneumonise - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
7	21	17155	16229	gnl P101623514	purative Pab protein (Bacillus subtilis)	62	46	927
7	24	19526	18519	gi 1276634	beta-ketacyl-ACP synthase III (Cuphea wrightii)	62	37	1008
12	7	5904	4702	gi 1573768	ATG-specific adenine glyoxylase (H. influenzae)	62	43	1203
12	9	8032	8793	gi 1591587	panochetate metabolism flavoprotein (Methanococcus jannaschii)	62	33	762
15	11	9678	9328	gi 1211513011	hypothetical 20.3% identical to the amino acid sequence IS1111 - Agrobacterium tumefaciens (strain R02) plasmid 71	62	43	351
17	5	2609	2442	gi 1551091	fr. Jannaschii predicted coding region M0374 (Methanococcus jannaschii)	62	43	168
17	5	3053	2835	gi 149570	role in the expression of lactacin F, part of the laf operon (Lactobacillus sp.)	62	44	219
22	10	8627	9318	gnl P10160580	similar to S. subtilis DnaI (Bacillus subtilis)	62	43	912
30	3	865	2043	gi 2314379	(AE000627) ABC transporter, ATP-binding protein (yhcG) (Haemobacter pylori)	62	43	1179
33	5	2235	1636	gi 113976	lipa-52r gene product (Bacillus subtilis)	62	44	600
38	11	5889	6123	gi 146231	0251 (Escherichia coli)	62	34	435
40	17	14272	13328	gnl P10161994	hypothetical protein (Methanococcus jannaschii)	62	43	945
42	1	3	311	gi 1146182	putative (Bacillus subtilis)	62	41	309
44	2	1287	6005	gi 1784952	(AE000176) 0877, 100% identical to the first 60 residues of the 100 aa hypothetical protein fragment YG04_E00145 SW: P45146	62	43	2739
46	12	9732	9304	gi 166920	repressor protein (Enterococcus hirae)	62	32	429
51	8	5664	7181	gnl P10161153	StyK1 methylase (Salmonella enterica)	62	44	1518
52	3	2791	2099	gi 1182686	integral membrane protein (Bacillus subtilis)	62	41	693
55	16	15702	14704	gnl P101613028	hypothetical protein (Bacillus subtilis)	62	40	999
59	6	3418	3984	gi 1505463	unknown (Lactococcus lactis lactis)	62	32	567
63	5	4897	4609	gi 149773	polin gene inverting protein (PivM) (Morganella morganella)	62	28	189
70	14	10002	10729	gi 192977	bplC gene product (Bordetella pertussis)	62	45	738
71	13	18790	20383	gi 1380135	located for by C. elegans cDNA library; similar to melibiose carrier protein (thiogenin) (Agrobacterium tumefaciens)	62	62	1593
71	28	32217	32768	gnl P10161312	Ygac (Bacillus subtilis)	62	35	552
74	7	11666	10983	gi 1552753	hypothetical (Escherichia coli)	62	38	1286

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% hom	Match (nt)
90	8	9310	9509	gi 121010002	[AB01188] FUNCTION UNKNOWN. [Bacillus subtilis]	62	46	240
97	10	9068	7001	gi 1842163	protein-N(p)-phosphatidylase sugar phosphatase [Escherichia coli]	62	42	208
98	4	2306	3268	gi 121010496	IsaE (integral membrane protein) [Pseudomonas aeruginosa]	62	42	963
102	3	2823	3519	gi 121010100	hypothetical protein [Bacillus subtilis]	62	24	711
103	3	2795	1242	gi 121010049	H. influenzae hypothetical ABC transporter; P4808 (974) [Bacillus subtilis]	62	41	1354
111	2	2035	3462	gi 581297	[NisP [Lactococcus lactis]	62	44	1428
112	6	4939	5649	gi 1574381	[IleC-1 operon protein (IleC) [Haemophilus influenzae]	62	39	923
112	6	4939	5649	gi 1574381	[IleC-1 operon protein (IleC) [Haemophilus influenzae]	62	39	923
124	3	1137	721	gi 1573024	energetic ribonucleoside-triphosphate reductase (ncd) [Haemophilus influenzae]	62	45	417
124	6	3162	2329	gi 609076	[lucyl aminopeptidase [Lactobacillus delbrueckii]	62	40	834
126	7	11073	7516	gi 121010163	[ORF4 [Bacillus subtilis]	62	38	1558
129	6	4983	4540	gi 1641509	[zinc finger protein Efe - Chilo iridescent virus	62	48	444
131	7	4510	4103	gi 1457245	[Lactococcus lactis]	62	42	408
149	2	1923	2579	gi 1592142	ABC transporter, probable ATP-binding subunit [Methanococcus jannaschii]	62	41	657
149	7	5340	6055	gi 1210102500	[YidO protein [Bacillus subtilis]	62	40	996
156	1	450	238	gi 12101025644	[membrane protein [Streptococcus pneumoniae]	62	40	213
159	6	3606	2935	gi 1210102050	[transmembrane [Bacillus subtilis]	62	37	672
171	2	1779	2291	gi 43941	[EII-B for PTS [Klebsiella pneumoniae]	62	35	513
172	2	385	723	gi 1895550	[putative cellobiose phosphotransferase enzyme III [Bacillus subtilis]	62	39	339
173	3	2599	493	gi 1591732	[cobalt transport ATP-binding protein O [Methanococcus jannaschii]	62	42	1107
179	2	492	1754	gi 1574071	[H. influenzae predicted coding region H1038 [Haemophilus influenzae]	62	38	1163
181	6	2856	3707	gi 1777435	[LacR [Lactococcus casei]	62	42	852
185	2	2074	311	gi 2182397	[AE000073] Y41N [Rhizobium sp. NGR234]	62	41	1764
200	2	1061	1984	gi 1450566	[transmembrane protein [Bacillus subtilis]	62	37	924
202	3	2543	3473	gi 422219	[P35 gene product (AA 1 - 314) [Escherichia coli]	62	41	891
210	3	1374	1565	gi 49315	[ORF1 gene product [Bacillus subtilis]	62	45	192

TABLE 2  
S. pneumoniae - putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
211	1	3	971	gi147402	hemolysin 117-kDa (Escherichia coli)	62	43	969
223	2	1495	1034	gnl pid d101190	ORF2 (Streptococcus mutans)	62	41	462
228	1	34	90	gi1510063	glycerol uptake facilitator (Streptococcus pneumoniae)	62	44	876
234	2	90	917	gi1293259	(AF008220) YqkI (Bacillus subtilis)	62	38	828
242	5	1765	1497	gi1293259	galactosyltransferase (Arabidopsis thaliana)	62	33	279
375	1	1	159	gi11674231	(AE000092) Mycoplasma pneumoniae, hypothetical protein homolog; similar to Seta-Prot Accession Number P35155, from B. subtilis (Mycoplasma pneumoniae)	62	40	159
385	5	584	357	gi11573353	outer membrane integrity protein (toxA) (Haemophilus influenzae)	62	47	228
39	14550	3269	gi1606162	ORF_229 (Escherichia coli)		61	41	720
7	4	2725	3225	gi12114425	similar to Synchocystis sp. hypothetical protein, encoded by Genbank Accession Number D64006 (Bacillus subtilis)	61	43	501
17	6	3326	1054	gi1149549	Lactacin P (Lactobacillus sp.)	61	43	273
44	3	4061	4957	gnl pid d101068	xylose repressor (Synchocystis sp.)	61	38	897
54	11	8398	7234	gnl pid d101326	YqkH (Bacillus subtilis)	61	42	1155
57	6	3974	6037	gnl pid d101116	YqkI (Bacillus subtilis)	61	42	2044
58	5	7355	5545	gi1245169	SPYDING/POWESKINE TRANSPORT SYSTEM PROTEIN POC.	61	34	792
67	1	3	692	gi1537108	ORF_254 (Escherichia coli)	61	46	690
68	9	8816	7890	gi119503	IP94212 gene product (AA 1-144) (Lupinus polyphylus)	61	41	927
70	15	10737	12098	gi1992976	bpIF gene product (Bordetella pertussis)	61	44	1272
72	11	9759	12022	gnl pid d101833	carboxymagnaporidine decarboxylase (Synchocystis sp.)	61	36	444
76	8	7881	7053	gnl pid d100305	(arney) diphosphate synthase (Bacillus stearothermophilus)	61	45	879
87	4	4914	3697	gi1528991	unknown (Bacillus subtilis)	61	42	1218
87	13	12311	11361	gi11769681	(AE000407) methionyl-tRNA formyltransferase (Escherichia coli)	61	44	951
91	2	731	3499	gi1537080	ribonucleoside triphosphate reductase (Escherichia coli)	61	45	2259
105	3	2711	3499	gnl pid d101851	hypothetical protein (Synchocystis sp.)	61	44	789
115	6	7966	4478	gi1892747	putative cel operon regulator (Bacillus subtilis)	61	36	1491
123	8	7181	9518	gi12099227	protein bisecting kinase (Bacteroides fecalis)	61	40	1338

TABLE 2

S. pneumoniae - putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
126	4	7525	6725	gi 1387043	(AF000184) 271; This 271 aa orf is 24 pct identical (16 aa) to 245 residues of an approx. 272 aa protein YIDA_ECOLI SM: P09397 [Escherichia coli]	61	38	801
128	1	1	639	gn P12D d101328	Y04Y [Bacillus subtilis]	61	41	639
139	7	4794	5034	gi 1022726	unknown [Staphylococcus haemolyticus]	61	41	261
139	9	12632	5913	gn P12D a720014	beta-1,3-galactosyl transferase [Bacteroides adhaerens]	61	41	6720
143	1	2552	42	gi 520541	penicillin-binding proteins 1A and 1B [Bacillus subtilis]	61	42	2511
146	16	13125	11424	gi 1532743	phosphoglucomutase [Bacillus subtilis]	61	42	702
162	3	4112	3456	gn P12D d101829	phosphoglucomutase [Synchocystis sp.]	61	30	657
172	3	727	1077	gn P12D d102048	phosphoglucomutase [Bacillus subtilis]	61	44	351
177	3	1101	1772	gn P12D d100374	unknown [Bacillus subtilis]	61	43	672
202	2	1278	2585	gi 1045831	hypothetical protein (GB1318945.6) [Myoplasma genitalium]	61	36	1108
224	3	2742	3144	gi 1591144	M. jamaeensis predicted coding region M0440 [Methanococcus jamaeensis]	61	30	363
225	4	3195	3766	gi 1532774	hypothetical [Escherichia coli]	61	40	572
249	2	212	802	gi 1000453	Trar [Bacillus subtilis]	61	42	591
254	2	843	484	gn P12D d100417	ORF120 [Escherichia coli]	61	36	360
257	1	3	350	gn P12D e255315	unknown [Mycobacterium tuberculosis]	61	42	348
293	4	3971	3657	gi 1251513 C61	hypothetical 30.3K aa protein (human sequence 191131) - Agrobacterium tumefaciens (strain 9023) plasmid p1	61	45	315
301	1	949	17	gi 22931209	(AF016424) contains similarity to acyltransferases [Caenorhabditis elegans]	61	33	933
373	1	1066	287	gi 1393396	TP-232 membrane associated protein [Trypanosoma brucei sub-group]	61	38	780
3	16	12473	124955	gi 1537093	ORF_0153b [Escherichia coli]	60	27	483
6	5	4636	5739	gi 2293258	(AF008220) YcoI [Bacillus subtilis]	60	35	1104
6	12	119336	111867	gi 2930107	ORF3 (put.); putative [Laetococcus lactis]	60	44	750
17	13	6708	6484	gi 149569	Lactacin F [Lactobacillus sp.]	60	32	255
18	7	6977	5670	gi 1788140	(AF000278) 0481; This 481 aa orf is 35 pct identical (19 aa) to 309 residues of an approx. 856 aa protein B041_HUMAN SM: P46087 [Escherichia coli]	60	43	1108
20	15	15878	17167	gn P12D d100584	unknown [Bacillus subtilis]	60	44	1290

TABLE 2

5. pneumoniae Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match gene name	% sim	% ident	Length (nt)
22	1	1	243	gi 12102050	transmembrane [Bacillus subtilis]	60	36	243
32	10	8286	8964	gi 2932375	[AF008220] YtaG [Bacillus subtilis]	60	37	669
38	15	8837	9497	gi 40023	ip subtilis genes rpmk, rpa, 50kd, glid and glod [Bacillus subtilis]	60	35	861
43	6	8610	9544	gi 171787	protein kinase 1 [Saccharomyces cerevisiae]	60	36	2667
44	1	1	1269	gi 12102050	unknown [Saccharomyces cerevisiae]	60	44	1269
45	10	11138	10368	gi 397488	1,4-alpha-glucan branching enzyme [Bacillus subtilis]	60	43	771
48	19	15766	14378	gi 121020173	orf1 [Bacillus subtilis]	60	39	1389
48	21	16727	16551	gi 121020141	unannotated protein product [Haemophilus actinomycetemcomitans]	60	32	225
50	2	898	1177	gi 121020173	orf1 [Bacillus subtilis]	60	39	1389
62	2	638	1177	gi 121020173	orf1 [Bacillus subtilis]	60	42	540
68	350	5203	5203	gi 121020173	orf1 [Bacillus subtilis]	60	42	540
70	11	5781	6182	gi 12102014	[AB001088] SIMILAR TO YDFR GENE PRODUCT OF THIS ENTRY (YDFR_BACSDU)	60	33	402
70	12	6143	8133	gi 121020173	hypothetical protein [Bacillus subtilis]	60	38	1791
71	8	11701	14157	gi 580866	ipa-126 gene product [Bacillus subtilis]	60	33	2457
74	8	12509	11664	gi 1210201832	phosphatidase cytidyltransferase [Synecocystis sp.]	60	45	846
76	4	4116	3357	gi 2352096	orf1; similar to serine/threonine protein phosphatase [Fervidobacterium islandicum]	60	39	750
80	4	7372	7665	gi 1784420	[AB000131] f66; 100 pct identical to GB: EC00H21.6 ACCESSION: D18582 [Escherichia coli]	60	30	294
81	6	4073	4522	gi 1147402	hypothetical protein [Bacillus subtilis]	60	35	450
86	1	940	155	gi 143177	putative [Bacillus subtilis]	60	26	786
92	1	1	1	gi 1784389	hypothetical protein [Escherichia coli]	60	45	192
93	14	10619	9184	gi 1784389	[AB000297] o464; This 464 aa orf is 33 pct identical (9 gaps) to 331 residues of an approx. 416 aa protein HPRC-NE100 SH: P43905 [Escherichia coli]	60	27	1236
94	5	5548	8121	gi 121020985	cyclic nucleotide-gated channel beta subunit [Rattus norvegicus]	60	50	2574
97	7	5396	4533	gi 1591396	transketolase [Methanococcus jannaschii]	60	43	864
102	2	2081	2813	gi 121020929	hypothetical protein [Mycobacterium tuberculosis]	60	43	753

TABLE 2

S. pneumoniae - Relative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match	Match accession	Match gene name	% sim	% ident	length (nt)
106	9	9773	9183	gnl P1D e334782	y18M protein [Bacillus subtilis]		60	31	593
113	8	6361	6837	gi 1466575	inf10; inf96; CL137 [Mycobacterium leprae]		60	43	477
115	2	2755	534	gnl P1D e328143	(A000332) glucosidase II [Homo sapiens]		60	32	232
122	7	4753	5068	gnl P1D d101876	transposase [Synecocystis sp.]		60	39	306
127	8	4510	5293	gi 1777938	Pgm [Treponema pallidum]		60	38	774
138	4	3082	2672	gnl P1D e325196	hypothetical protein [Bacillus subtilis]		60	36	411
139	1	177	4	gnl P1D d100680	ORF [Thermus thermophilus]		60	39	174
139	11	14520	13069	gi 1537145	ORF 4437 [Escherichia coli]		60	30	1312
140	2	2592	1249	gi 1205627	protein histidine kinase [Enterococcus faecalis]		60	37	1344
141	1	210	1049	gi 1463181	B5 ORF from bp 3842 to 4081; putative [Human papillomavirus type 33]		60	34	840
141	5	5368	6005	gi 145162	tyrosine-sensitive DAP synthase [arop] [Escherichia coli]		60	41	1030
142	6	3558	4049	gi 1600711	putative [Bacillus subtilis]		60	37	492
148	10	7742	8713	gnl P1D e110022	hypothetical protein [Bacillus subtilis]		60	27	972
153	5	3667	4278	gi 1293322	(AF098220) branch-chain amino acid transporter [Bacillus subtilis]		60	42	612
155	1	1413	748	gi 12104504	putative UDP-glucose dehydrogenase [Escherichia coli]		60	40	686
158	3	3116	2472	gnl P1D d100872	a negative regulator of pho regulon [Pseudomonas aeruginosa]		60	37	645
159	3	776	1386	gnl P1D e300090	phoA, highly similar to Bacillus anthracis CapA protein [Bacillus subtilis]		60	48	609
163	7	8049	8668	gnl P1D d103313	YqjH [Bacillus subtilis]		60	38	420
170	3	4130	2688	gi 1574179	flr influenzae predicted coding region H1244 [Haemophilus influenzae]		60	39	1443
171	7	4717	5201	gi 1606076	ORF_0384 [Escherichia coli]		60	44	1185
183	3	2440	2195	gi 1877427	repressor [Streptococcus pyogenes phage T12]		60	38	306
191	10	9444	8828	gi 1415664	catabolite control protein [Bacillus megaterium]		60	42	1017
200	1	139	1093	gi 148462	transmembrane protein [Bacillus subtilis]		60	37	945
201	3	3895	1928	gi 1475112	enzyme Iiabc [Pediococcus pentosaceus]		60	39	1968
214	15	11990	10439	gi 1877407	hypothetical [Haemophilus influenzae]		60	39	492
218	4	2145	2363	gi 1608520	myosin heavy chain kinase A [Dictyostellium discoideum]		60	31	219



TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length
218	4	2518	2351	gi 437705	hyaluronidase (Streptococcus pneumoniae)	60	53	168
242	1	725	3	gi 43938	Sor regulator (Klebsiella pneumoniae)	60	41	72
245	1	1	288	gi 304897	lecA type I restriction modification enzyme M subunit (Escherichia coli)	60	56	288
251	1	905	45	gi 61632	unknown (Staphylococcus aureus)	60	35	85
259	1	869	82	gi 133794	reg (Streptococcus gordonii)	60	32	888
260	2	1482	1662	gi 531840 E318	probable transposase - Bacillus stearothermophilus	60	36	171
274	1	836	96	gi 1592173	W-ethylamine chlorohydrolase Methanococcus jannaschii	60	40	741
308	1	463	2	gi 1787397	(AF000214) o157 Escherichia coli	60	43	462
318	1	3	308	gi 1293147	kerc recombinase (Lactobacillus leichmannii)	60	42	306
344	1	73	522	gi 509872	repressor protein (Bacteriophage Tn-2009)	60	32	450
5	1	576	6	gi 2933147	(AF008220) TxAH Bacillus subtilis	59	31	573
7	22	18140	17142	gi 1210 e280724	unknown (Mycobacterium tuberculosis)	59	39	993
10	1	3413	6	gi 1353880	isialidase L (Mecrobactria decora)	59	41	1410
15	6	6463	5156	gi 560941	PII Bacillus subtilis	59	35	1308
22	2	479	1393	gi 142469	els operon regulatory protein (Bacillus subtilis)	59	34	915
22	5	2698	4614	gi 1210 e280423	PCPA (Streptococcus pneumoniae)	59	44	1317
30	1	208	558	gi 1210 e23868	hypothetical protein (Bacillus subtilis)	59	37	351
30	4	3678	2455	gi 1210 e202390	unknown (Lactobacillus sake)	59	33	1224
35	13	12201	11071	gi 1210 e238664	hypothetical protein (Bacillus subtilis)	59	35	1131
35	14	13248	12183	gi 1051647	CusH (Staphylococcus aureus)	59	39	1107
36	18	18076	17897	gi 1500535	M. jannaschii predicted coding region W1635 Methanococcus jannaschii	59	33	180
38	12	8172	7137	gi 2293239	(AF008220) TxAH Bacillus subtilis	59	34	966
42	3	1952	3361	gi 1684845	pinin (Canis familiaris)	59	40	1410
50	3	4878	1778	gi 1210 d01329	VqJK (Bacillus subtilis)	59	41	951
56	5	1870	2188	gi 1210 e17594	kerc recombinase (Lactobacillus leichmannii)	59	29	519
61	6	6812	5648	gi 1210 e11516	aminotransferase (Bacillus subtilis)	59	40	1185
67	5	2382	3023	gi 1146190	2-keto-3-deoxy-6-phosphogluconate aldolase (Bacillus subtilis)	59	36	642

TABLE 2  
S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	Length (nt)
69	10	8547	8849	gi1572628	Leucine binding (Coat) [Haemophilus influenzae]	59	38	253
87	12	11383	10055	gn1 P1D e023504	Putative Fmu protein [Bacillus subtilis]	59	44	1329
113	14	13927	13894	gi1473731	(A800010) Mycoplasma fermentans, similar to Splice Prot. Accession Number P20946, from E. coli [Mycoplasma pneumoniae]	59	43	1368
115	15	8531	8531	gi1586086	R. Jannaschii predicted coding region M0110 [Methanococcus Jannaschii]	59	38	246
119	2	1946	1526	gn1 P1D e030005	homologue to ORF2 in nrfEF operons of E. coli and S. typhimurium [Lactococcus lactis]	59	43	441
138	17	11438	13178	gn1 P2D e079632	unknown [Mycobacterium tuberculosis]	59	38	261
140	22	23903	23388	gi1482922	protein with homology to palI repressor of B. subtilis [Lactobacillus delbrueckii]	59	40	516
148	13	9657	9014	gn1 P1D d102005	(A800148) FUNCTION UNKNOWN, SIMILAR PRODUCT IN H. INFLUENZAE AND STREPTOCOCCUS [Bacillus subtilis]	59	32	684
149	10	7213	8244	gi1710422	cap-binding factor 1 [Staphylococcus aureus]	59	40	1032
164	9	6593	6013	gn1 P1D d100965	ferric anguibactin-binding protein precursor FABC of V. anguillarum [Bacillus subtilis]	59	41	981
164	12	8836	7823	gn1 P1D d100954	homologue of ferric anguibactin transport system permease protein PABC of V. anguillarum [Bacillus subtilis]	59	35	1014
177	2	401	1072	gi1289739	Protein for C. pneumoniae (Genbank Z14728); putative [Carnobacterium elegans]	59	40	672
177	7	3841	4200	gi12313445	(A8000551) H. pylori predicted coding region HP0342 [Helicobacter pylori]	59	38	360
183	6	2768	2508	gi1509672	repressor protein [Bacteriophage T4c009]	59	50	261
186	6	3398	2820	gi1660680	ORF_2590; Geneplot suggests frameshift linking to o267, not found [Escherichia coli]	59	38	579
190	3	3120	1711	gi1613768	histidine protein kinase [Streptococcus pneumoniae]	59	32	1410
194	2	1621	1019	gn1 P1D d100579	unknown [Bacillus subtilis]	59	40	603
198	7	5205	4106	gn1 P1D e013073	hypothetical protein [Bacillus subtilis]	59	38	900
220	5	4362	3958	gn1 P1D d101322	YqjH [Bacillus subtilis]	59	46	405
242	3	1573	2167	gi1187045	(A8000184) f308; This 308 aaorf is 35 pnt identical (35 aa) to 305 residues of an approx. 296 aa protein PFAC_EXULI SM: P12675 [Escherichia coli]	59	42	795
247	2	1154	1460	gi140073	gn1 P1D d10073	59	39	327

TABLE 2  
S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
256	1	468	2	gnl PTD d101924	hemolysin [Synychochysis sp.]	59	39	867
258	1	65	420	gi 1246532	ORF 73, contains large complex repeat CR 73 [Kaposi's sarcoma associated herpesvirus]	59	20	756
270	1	386	1126	gnl PTD d102092	YnfB [Bacillus subtilis]	59	40	741
281	1	532	166	gi 166062	putative [aerococcus lactis]	59	31	387
309	1	3	479	gi 1405879	YnfH [Escherichia coli]	59	38	477
363	1	2	1894	gi 915208	gastric mucin [Sus scrofa]	59	31	1893
387	2	425	84	gi 160671	S antigen precursor [Plasmodium falciparum]	59	44	342
5	6	11223	10465	gnl PTD d101812	LuoQ [Synychochysis sp.]	58	29	759
29	4	2098	3513	gnl PTD d100479	Rax - ATPase subunit J [Enterococcus hirae]	58	39	1416
30	5	4058	3651	gi 139478	ATP binding protein of transport ATPases [Bacillus firmus]	58	34	408
33	6	2983	2210	gnl PTD d101164	unknown [Bacillus subtilis]	58	45	774
36	8	5316	6179	gi 1518679	orf [Bacillus subtilis]	58	32	864
43	5	5926	3971	gi 1788150	[Aeromonas aeruginosa]	58	37	3956
46	1	3704	5221	gnl PTD e267229	unknown [Bacillus subtilis]	58	42	1518
48	14	11722	11066	gnl PTD d101771	thiamin biosynthetic bifunctional enzyme [Synychochysis sp.]	58	34	657
52	1	1229	3	gnl PTD d101291	reductase [Pseudomonas aeruginosa]	58	35	1227
53	2	702	412	gi 1313357	cytochrome c biogenesis protein (ccbA) [Halobacterium pyroclit]	58	25	291
58	4	6586	5458	gi 1247129	transport protein [Escherichia coli]	58	41	1089
69	5	4914	3807	gnl PTD d101192	unknown [Bacillus subtilis]	58	41	1128
71	27	31357	32277	gi 2408014	hypothetical protein [Schistosoma haematophyllum]	58	33	921
72	4	3586	3882	gi 116094	modulin [Escherichia coli]	58	34	705
74	3	4937	4210	gi 12293252	[Aeromonas hydrophila]	58	33	708
79	4	4594	1422	gnl 1217999	ORF3 [Streptococcus pneumoniae]	58	44	1173
82	8	10585	8171	gi 1882711	ironuclease V alpha subunit [Escherichia coli]	58	38	2415
84	17	18602	12317	gi 147442	3-dehydroquinate hydrolyase (3-dehydroquinate)	58	32	681
97	2	931	560	gi 1513794	rdp [Streptococcus gordonii]	58	32	372

TABLE 2

S. pneumoniae - putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Search accession	Match gene name	% sim	% ident	Length
108	2	358	2724	gi1537020	vecB gene product [Escherichia coli]	58	37	2367
111	5	4593	5240	gi1592142	ABC transporter, probable ATP-binding subunit [Methanococcus jannaschii]	58	35	648
120	3	4421	5110	gn1 PFI d101330	YogX [Bacillus subtilis]	58	47	690
128	16	13131	12673	gi1662919	OMP U [Enterococcus hirae]	58	42	459
132	3	6174	4939	gi18800101	macrolide-efflux determinant [Streptococcus pneumoniae]	58	35	1136
133	1	111	890	gn1 PFI d0269488	Unknown [Bacillus subtilis]	58	36	780
160	11	8615	9865	gi1473901	ORF1 [Lactococcus lactis]	58	39	1351
161	6	6248	6849	gn1 PFI d101024	ED-1 protein [Homo sapiens]	58	32	582
169	1	214	2	gn1 PFI d100447	translation elongation factor-3 [Chlorella virus]	58	31	213
187	1	487	2	gi1475114	regulatory protein [Methanococcus putrescentiae]	58	38	486
187	6	4384	4620	gi1167475	desiccation-related protein [Crotosigma plantagineum]	58	55	237
190	2	1464	1640	gn1 PFI d146272	serpentine pheromone [Streptococcus gordonii]	58	38	177
192	2	2012	1344	gn1 PFI d100856	[rat OCT360 [Rattus rattus]]	58	44	669
206	1	1292	69	gn1 PFI d020977	product similar to WRB [Lactobacillus sakei]	58	35	597
216	2	3333	555	gn1 PFI d025036	hypothetical protein [Bacillus subtilis]	58	33	1779
217	5	3220	4821	gi1468474	cellulose phosphoryltransferase enzyme II** [Bacillus stearothermophilus]	58	38	930
217	7	5636	5106	gn1 PFI d102048	B. subtilis cellulose phosphoryltransferase system cell: P46317 (994)	58	44	531
232	1	2	811	gi1573777	cell division ATP-binding protein (fap) [Methanophila influenzae]	58	39	810
264	1	2	715	gi1973330	NATA [Bacillus subtilis]	58	32	714
280	1	33	767	gi1186187	hypothetical 29.6 kD protein in DMC-talB intergenic region [Escherichia coli]	58	31	735
306	1	845	3	gn1 PFI d0334760	YibC protein [Bacillus subtilis]	58	47	843
360	3	1556	1092	spP46351 Y2CD	HYPOHYPHETAL 45.4 kD PROTEIN IN THIAMINASE 1 5'-REGION	58	32	465
363	5	2160	1867	gi1160671	S antigen precursor [Plasmodium falciparum]	58	51	294
372	1	806	3	gi1393394	Tb-231 membrane associated protein [Trypanosoma brucei]	58	37	804
382	2	749	519	gi1 PFI d151 JC11	hypothetical 20.3K protein (insertion sequence [S131]) - Agrobacterium tumefaciens (strain 1022) plasmid T1	58	41	231

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	Length (nt)
3	9	8609	7471	gi 1193745	M. j. j. predicted coding region M0912 (Methanococcus jannaschii)	57	38	935
10	10	7634	7507	gi 1737169	homologue to SRF (Arabidopsis thaliana)	57	30	168
1	2	412	412	gn 12100139	ORF (Methanobacter pasteurianus)	57	42	411
31	4	2032	1388	gi 2493213	(AF088220) YrpA (Bacillus subtilis)	57	37	645
31	5	6549	6549	gn 12100139	hypothetical protein [Bacillus subtilis]	57	36	483
45	5	5446	5060	gi 1592204	phosphatase phosphatase (Methanococcus jannaschii)	57	44	387
2	2	635	752	gi 155369	PTS enzyme-II fructose 1,6-bisphosphate	57	35	1110
52	6	4520	6850	gi 1574144	single-stranded-DNA-specific exonuclease (recJ) (Methanophilus influenzae)	57	35	2331
51	5	2079	1795	gi 1643580	replicase-associated polypeptide (coat blue dwarf virus)	57	46	285
63	6	5312	4995	gi 2182608	(AE000941) YnfJ (Rhizobium sp. NGR234)	57	39	318
72	15	13883	13059	gn 121010892	homologue to BuisProt.YnfA, BCO11 hypothetical protein [Bacillus subtilis]	57	40	825
79	2	2561	1815	gn 121010965	homologue of NADH-flavin oxidoreductase Fp of V. harvey [Bacillus subtilis]	57	44	747
82	9	9596	9763	gi 1206045	short region of similarity to glycerophosphoryl diester phosphodiesterases (Caenorhabditis elegans)	57	35	168
86	16	15371	14493	gi 1787983	(AE000264) c288; 92 pct identical (1 gap) to 222 residues of fragment (125 aa) [Escherichia coli]	57	34	879
93	3	1695	1177	gi 1500003	initiator mut protein (Methanococcus jannaschii)	57	33	519
96	6	3026	4519	gi 1593682	threonine synthase [Arabidopsis thaliana]	57	43	1494
99	14	17211	18212	gi 1773349	BlaA protein [Bacillus subtilis]	57	44	1002
112	6	7448	7903	gi 1593393	M. jannaschii predicted coding region M0978 (Methanococcus jannaschii)	57	30	456
113	16	18627	18328	plf145605/1456	mature-parasite-infected erythrocyte surface antigen ME2A - Plasmodium falciparum	57	22	300
123	2	343	1110	gi 168141754	hypothetical protein M0355 (Methanophilus influenzae strain Rd M20)	57	38	768
123	4	2108	2884	gn 1210102148	(AE001684) sulfate transport system permease protein (Chlorella vulgaris)	57	39	777
127	15	6577	5597	gi 1575062	crotonase C (mtc) (Methanophilus influenzae)	57	35	891
128	13	9251	9790	gi 153692	pneumolysin (Streptococcus pneumoniae)	57	38	540
131	3	3319	1363	gi 42081	hgd gene product (AA 1-250) [Escherichia coli]	57	36	777

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

TABLE 2

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
136	1	214	1221	[bb]148453	Spaa-endocytosis membrane protein [Streptococcus subinus, MUCOB 25, Peptide, 1566 aa] [Streptococcus subinus]	57	44	1008
140	25	28701	28851	[g]1505576	beta-glucoside permease [Bacillus subtilis]	57	38	1851
141	6	6395	7418	[g]195560	unknown [Schizosaccharomyces pombe]	57	44	1044
144	3	3231	2785	[gm]PRD100139	[ORF] [Aerobacter pasteurianus]	57	42	447
155	4	5454	4554	[g]1600431	glyoxyl transferase [Erebia amygdalora]	57	34	591
159	9	4877	5854	[g]1506509	[g]307 [Escherichia coli]	57	35	978
167	11	9710	9249	[gm]PRD100139	[ORF] [Aerobacter pasteurianus]	57	42	482
171	6	4023	4436	[g]147402	hamase permease subunit III-Man [Escherichia coli]	57	29	414
178	4	2170	1076	[gm]PRD102004	[AB001488] ATP-DEPENDENT RNA HELICASE RPD-102000. [Bacillus subtilis]	57	39	1095
190	1	145	1455	[g]149420	export/processing protein [Lactococcus lactis]	57	30	1311
198	1	298	95	[g]1522268	unidentified ORF22 [Bacteriophage b1607]	57	36	204
203	2	3195	2110	[gm]PRD1283915	orf c01003 [Sulfolobus solfataricus]	57	41	1086
205	1	40	507	[g]1439527	ETA-man [Lactobacillus acidophilus]	57	28	468
214	7	4243	3797	[gm]PRD102049	[H. influenzae, ribosomal protein alanine acetyltransferase; P4305 (189) [Bacillus subtilis]	57	48	417
268	3	1767	1276	[g]143979	L-curvatus small cryptic plasmid gene for rep protein [Lactobacillus curvatus]	57	36	492
351	1	324	34	[gm]PRD1274871	FO386.b [Comoriabactis elegans]	57	31	291
386	1	226	2	[g]160671	[S antigen precursor [Plasmodium falciparum]	57	45	225
5	5	10486	877	[g]105895	YnfH [Escherichia coli]	56	33	1710
8	5	3674	3910	[g]1467199	pKac; L518.P.2 Mycobacterium leprae]	56	39	231
10	2	3442	1874	[gm]PRD101907	sodium-coupled permease [Synchytrium sp.]	56	36	1569
21	1	1880	333	[g]1211949	[AE005931] oncoprotection protein (pomp2) [Helicobacter pylori]	56	33	1548
22	29	14968	22456	[gm]PRD102001	[AB001488] FIBRILLAR ACETYLTRANSFERASE. [Bacillus subtilis]	56	37	489
27	1	1361	3	[g]1215132	ecp59 (525) [Bacteriophage lambda]	56	30	1559
28	9	4687	4278	[g]1591090	DNA repair protein RAD2 [Methanococcus jannaschii]	56	29	390
33	1	3	386	[gm]PRD100139	[ORF] [Aerobacter pasteurianus]	56	41	384

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match	Gene name	% sim	% ident	length
36	7	5122	5337	[gi 170053 7000		hypothetical protein (proc 3' region) - <i>Pseudomonas aeruginosa</i> (strain PMO) (fragment)	56	28	276
40	4	3137	4318	[gi 1860301		macrolide-efflux determinant [ <i>Streptococcus pneumoniae</i> ]	56	27	1182
40	16	12511	13191	[gn P12 d127602		[P12] <i>Lactobacillus plantarum</i>	56	38	681
48	17	13775	13021	[gi 143729		transcription activator [ <i>Bacillus subtilis</i> ]	56	35	753
75	4	1674	2294	[gn P12 d102036		hypothetical protein [ <i>Bacillus stearothermophilus</i> ]	56	25	921
85	3	1842	1459	[gn P12 d100139		[ORF] [ <i>Neisseria pasteuriana</i> ]	56	41	364
89	7	5815	4940	[gi 163777		hypothetical protein to E. coli PPA protein [ <i>Bacillus subtilis</i> ]	56	42	876
105	2	1360	2718	[gn P12 d101913		hypothetical protein [ <i>Synechocystis</i> sp.]	56	37	1359
112	3	2151	3194	[gi 1527201		[ORF-645] [ <i>Escherichia coli</i> ]	56	31	1044
113	4	2754	2963	[gn P12 d100340		[ORF] [flum pox virus]	56	28	210
122	3	1203	2054	[gi 144035		hypothetical periplasmic glutamine binding protein [ <i>Salmonella typhimurium</i> ]	56	30	852
124	8	3939	3694	[gn P12 d248893		unknown [ <i>Mycobacterium tuberculosis</i> ]	56	27	246
125	4	4403	4107	[gn P12 d100247		human non-muscle myosin heavy chain (Myosin-10)	56	32	297
127	11	6608	6405	[gi 2182397		[AE000073] Y4H1 [ <i>Brizobium</i> sp. MGR214]	56	35	204
134	5	4769	3849	[gn P12 d101810		hypothetical protein [ <i>Synechocystis</i> sp.]	56	39	921
137	10	6814	7245	[gi 1592011		lipidate permease (cyt) [ <i>Methanococcus jannaschii</i> ]	56	34	412
142	8	3019	4582	[gn 151071 1420		hypothetical protein to E. coli, activity 5' of nifS - [ <i>Bacillus subtilis</i> ]	56	29	438
146	8	4676	3660	[gn P12 d101911		hypothetical protein [ <i>Synechocystis</i> sp.]	56	32	1017
148	3	1906	2739	[gn P12 d101099		phosphate transport system permease protein PexA [ <i>Synechocystis</i> sp.]	56	36	834
150	4	4449	2743	[gn P12 e104628		probably site-specific recombinase of the resolvase family of enzymes [bacteriophage TP21]	56	27	1707
172	1	2	238	[gi 1767791		[AE000249] F317, this 317 aa orf is 22 pct identical (11 aa) to 301 residues of an approx. 320 aa protein Y2C_BAC50 SW P39140 [ <i>Escherichia coli</i> ]	56	34	207
172	7	4979	5668	[gi 396293		[similar to <i>Bacillus subtilis</i> hypoch. 20 kDa protein, in tar 3' region [ <i>Escherichia coli</i> ]	56	40	690
186	7	3732	3197	[gi 1132200		PTS permease for mannose subunit rIPMan [ <i>Vibrio fischeri</i> ]	56	36	368
187	2	2402	819	[pk 1557904 5579		YnfR49 protein - <i>Streptococcus pyogenes</i> (strain CS101) (enocope M49)	56	35	1384

TABLE 2

S. pneumoniae - Relative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Search accession	Match gene name	% sim	% ident	length (aa)
204	3	2772	2239	gi 696376	DfE_0162 [Escherichia coli]	56	35	534
206	2	3142	1613	gi 559861	ICM [Pseudomonas]	56	38	1710
219	3	1689	1096	gi 1146137	putative [Bacillus subtilis]	56	27	554
230	2	409	1485	gi 6032482603	hypothetical protein 2 (nr 5' region) - Streptococcus mitis (strain DME175, serotype F)	56	40	1077
233	4	2910	3468	gi 1041785	riboptery protein [Pseudomonas]	56	24	339
273	2	1543	2724	gi 141089	jap protein [Bacillus subtilis]	56	32	1182
353	1	1	516	gi 1786952	hypothetical protein [Bacillus subtilis]	56	41	516
359	1	87	641	gi 1786952	(AE000176) s877, 100 pct identical to the first 86 residues of the 100 aa hypothetical protein fragment Y8B8_ECOLI SM: P54746 [Escherichia coli]	56	46	555
363	7	4482	4198	gi 1573352	outer membrane integrity protein (OIA) [Haemophilus influenzae]	56	38	285
376	1	2	508	gi P1D(e)325031	hypothetical protein [Bacillus subtilis]	56	33	507
38	4	1824	1618	gi P1D(e)316518	negative regulator of pilQ regulon [Pseudomonas aeruginosa]	55	31	660
28	4	1824	1618	gi P1D(e)316518	STAT protein [Dictyostelium discoideum]	55	40	207
29	6	4496	5041	gi 1088261	unknown protein [Arabidopsis sp.]	55	31	546
38	16	9695	10702	gi 580905	B. subtilis genes rpmA, rpmB, 50kD, gldA and gldB [Bacillus subtilis]	55	31	1008
49	5	5727	6182	gi 1786951	(AE000176) heat-responsive regulatory protein [Escherichia coli]	55	29	456
51	4	2381	2241	gi P1D(e)101293	Y8B8 [Bacillus subtilis]	55	42	161
52	9	9640	10866	gi 153016	DfE 419 protein [Staphylococcus aureus]	55	23	1227
53	4	1813	1349	gi 896042	DapF [Borrelia burgdorferi]	55	30	465
60	5	4794	5756	gi 1499876	magnesium and cobalt transport protein [Methanococcus jannaschii]	55	38	963
71	9	14176	15408	gi 1857120	glycyl transferase [Mycobacterium tuberculosis]	55	40	1233
75	6	3189	4229	gi P1D(e)208890	NAD alcohol dehydrogenase [Bacillus subtilis]	55	44	1041
108	10	10488	9820	gi P1D(e)324997	hypothetical protein [Bacillus subtilis]	55	36	699
113	12	13273	13037	gi P1D(e)311496	unknown [Bacillus subtilis]	55	34	765
113	13	13007	131945	gi 1573423	1-phospho-fructokinase (PFK) [Bacillus subtilis]	55	39	939
126	5	6764	5907	gi 1790131	(AE000446) Hypothetical 29.7 kb protein in bpaA-gybB intergenic region [Escherichia coli]	55	37	858



TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	Length (nt)
126	3	2719	862	gml PJ0101453	α-2-mannosidase [Bacillus licheniformis]	55	35	1818
138	3	2593	1610	gml J42833	ORF2 [Bacillus subtilis]	55	37	984
140	6	6716	5633	gml PJ0100964	homologous of hypothetical protein in a repaycin synthesis gene cluster of Bacillus subtilis	55	26	1284
147	3	3854	2136	gml J72330	dihydrolysoamide dehydrogenase [Clostridium magnum]	55	39	1719
147	110	10204	8921	gml PJ0171078	dihydroxycitrate dehydrogenase [Bacillus subtilis]	55	38	1284
148	5	3430	4119	gml J200572	peripheral membrane protein Y [Escherichia coli]	55	29	690
148	6	4171	4650	gml J695769	transposase [Bacillus subtilis]	55	37	480
149	14	12564	11650	gml PJ0101329	fgd [Bacillus subtilis]	55	32	915
156	3	1113	550	gml J314496	ORF06537 conserved hypothetical integral membrane protein [Haemophilus influenzae]	55	34	564
159	110	6625	5897	gml J290533	similar to E. coli ORF adjacent to suc operon; similar to gntR class of regulatory proteins [Escherichia coli]	55	29	729
164	3	1784	2332	gml PJ01255118	hypothetical protein [Bacillus subtilis]	55	37	549
164	5	2732	3631	gml J0248	put. resolvase Tnp I (M1 - 244) [Bacillus thuringiensis]	55	35	750
164	11	7428	7216	gml PJ01249407	unknown [Mycobacterium tuberculosis]	55	38	213
167	5	2860	3345	gml J53092	involved in protein secretion [Bacillus subtilis]	55	28	516
186	5	2880	2563	gml J05080	ORF 0290; Geneplot suggests frameshift linking to o267, not found [Escherichia coli]	55	35	318
189	8	4311	5396	gml PJ010183450	hypothetical tsg protein [Bacillus subtilis]	55	32	1086
192	5	3270	3079	gml J119504	vitellogenin convertase [Aedes aegypti]	55	38	192
195	2	2451	1384	gml J574693	transferrin, peptidoglycan synthesis (murG) [Haemophilus influenzae]	55	33	1071
198	4	3013	2471	gml PJ01013074	hypothetical protein [Bacillus subtilis]	55	29	543
214	3	371	744	gml PJ0101741	transposase [Synecocystis sp.]	55	33	372
219	2	1115	456	gml J288301	ORF2 gene product [Bacillus megaterium]	55	30	660
263	7	3742	3449	gml J18137	ORF-4 product [Chlamydomonas reinhardtii]	55	48	309
285	1	2	839	gml PJ0100974	unknown [Bacillus subtilis]	55	40	828
286	1	650	249	gml J398444	ORF (18 kDa) [Vibrio cholerae]	55	31	402
297	2	1239	1696	gml J50848	gptc [Porphyromonas gingivalis]	55	39	668

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins<sup>a</sup> similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match gene name	% sim	% ident	length
209	2	218	982	gi 1574463	hypothetical (Haemophilus influenzae)	55	35	765
328	2	646	234	gi 571500	prohibitin (Saccharomyces cerevisiae)	55	27	425
330	1	1340	474	gi 359397	isoX (Escherichia coli)	55	29	867
364	3	2538	1546	gi 391394	Tb-391 membrane associated protein (Trypanosoma brucei subgroup)	55	36	993
386	2	911	105	gi 160671	S antigen precursor (Plasmodium falciparum)	55	40	837
3	5	4604	3624	gi 2293176	(AF080220) signal transduction protein kinase (Bacillus subtilis)	54	26	981
5	11	7746	7246	gi 146245	putative (Bacillus subtilis)	54	26	981
38	34	16213	17937	gi 1480439	putative transcriptional regulator (Bacillus acetobutylicus)	54	27	1725
40	8	5076	4662	gi 359989	methionyl-tRNA synthetase (Bacillus stearothermophilus)	54	35	195
43	4	3380	2367	gnl pid a148611	ABC transporter (Lactobacillus helveticus)	54	25	1614
52	10	10844	12103	gi 1762962	FemA (Staphylococcus simulans)	54	29	1260
57	1	3	512	gi 558177	endo-1,4-beta-xylanase (Collimonas flini)	54	36	510
58	3	4749	4246	gnl pid H01237	hypothetical (Bacillus subtilis)	54	29	504
71	7	10684	11703	gi 510255	orf3 (Escherichia coli)	54	31	1020
71	120	27546	27737	gi 202543	serotonin receptor (Rattus norvegicus)	54	31	192
72	2	844	1098	gi 148613	srna gene product (Plasmod P.)	54	37	255
72	7	7438	6695	gi 136496	recombinase (Moraxella bovis)	54	38	744
74	10	14043	13465	gi 1209342	ORF 3 gene product (Staphylococcus aureus)	54	32	579
74	12	16493	15955	gi 2317798	naturase-related protein (Pseudomonas alcaligenes)	54	30	489
86	3	2877	2155	gi 164988	orf5.6 possibly encodes the O unit polymerase (Salmonella enterica)	54	34	723
89	5	4433	3921	gi 147211	phnO protein (Escherichia coli)	54	41	511
90	1	3	444	gi 1217798	naturase-related protein (Pseudomonas alcaligenes)	54	30	462
96	10	8098	8510	gnl pid H02015	(AB001488) SIMILAR TO SALMONELLA TYPHIMURUM SIXTY GENE REQUIRED FOR SURVIVAL IN MACROPHAGE. (Bacillus subtilis)	54	32	453
97	6	4662	3604	gi 1591394	trkA (Staphylococcus aureus)	54	30	1059
106	11	10406	12010	gi 1606266	ORF_0617 (Escherichia coli)	54	32	1405
147	8	8663	7404	gnl pid H01615	ORF_10_011977, similar to (SwissProt Accession Number P77340) (Escherichia coli)	54	35	1260

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
171	4	2477	3223	gi 1439528	E1IC-man (Lactobacillus curvatus)	54	36	747
174	2	2068	1787	gnl P1D100518	motor protein (Homo sapiens)	54	35	282
188	1	526	1168	gnl P1D1025052	unknown Mycobacterium tuberculosis	54	31	663
198	5	3582	2884	gnl P1D1011074	hypothetical protein [Bacillus subtilis]	54	33	699
207	1	1	1641	gnl P1D1010113	hypothetical protein [Synchococcus sp.]	54	24	1641
210	1	2	655	gnl 2232306	[AF008220] Ymp (Bacillus subtilis)	54	28	654
225	2	966	2357	gnl P1D1030194	g1196.1 (Caenorhabditis elegans)	54	39	1332
241	1	1681	347	gnl P1D101813	hypothetical protein [Synchococcus sp.]	54	26	1335
263	2	907	1395	gnl P1D101886	transposase [Synchococcus sp.]	54	30	489
263	6	3450	2977	gi 160671	S antigen precursor [Plasmodium falciparum]	54	47	474
277	3	2517	1363	gi 1166926	unknown protein [Streptococcus mutans]	54	30	1155
305	1	858	4	gi 2233198	[AF008220] Ymp (Bacillus subtilis)	54	37	750
325	1	19	768	gi 2182507	[AB000663] Y4IH (Rhizobium sp. NGR234)	54	32	309
332	2	898	590	gi 1591815	NDF-ribosylglycohydrolase (drac) [Methanococcus jannaschii]	54	49	240
365	4	240	479	gi 1530878	amino acid feature: N-glycosylation sites, aa 43, 46, 48, 51, 53, 56, 58, 61, 63, 65, 67, 100, 102, 105, 107, 165, 168; amino acid feature: Rod protein domain, aa 169.. 340; amino acid feature: globular protein domain	54		
7	35	19702	19493	gnl P1D1025511	hypothetical protein [Bacillus subtilis]	53	32	210
23	3	2497	2033	gnl P1D102015	[AB001488] SIMILAR TO SALMONELLA TYPHIMURION SLT GENE REQUIRED FOR SURVIVAL IN MICROPHAGE. [Bacillus subtilis]	53	25	465
29	11	9642	10121	gi 143331	alkaline phosphatase regulatory protein [Bacillus subtilis]	53	31	1080
33	3	1479	1009	pir S10655106	hypothetical protein X - Pyrococcus woesei (fragment)	53	33	471
36	6	4583	5134	gnl P1D1016029	unknown Mycobacterium tuberculosis	53	30	552
38	14	8521	8898	gi 1548094	homologous to E. coli rnpA [Bacillus subtilis]	53	30	378
52	7	7007	8686	gi 1377631	unknown [Bacillus subtilis]	53	29	1660
54	17	17555	19264	gi 166069	orf2 gene product [Lactobacillus leichmannii]	53	36	2010
56	1	1	681	gi 1592266	restriction modification system S subunit [Methanococcus jannaschii]	53	32	681

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	End (nt)	Match	Match position	Match gene name	% sim	% ident	Length (nt)
57	10	9411	8487	91	1788543	(AD000310) f351; Residues 1-121 are 100 pct identical to Y00L_ECOLI SW: P33944 (122 aa) and aa 152-131 are 100 pct identical to Y00L_ECOLI SW: P33943 (Escherichia coli)	53	31	945
61	1	429	4	91	1702467	80024.12 (Caenorhabditis elegans)	53	33	426
71	1	572	4	91	1703594	[Pb-29] membrane associated protein (Trypanosoma brucei subgroups)	53	33	5749
72	3	894	2840	91	1223178	(AF000220) YGAP (Bacillus subtilis)	53	27	1947
73	14	793	7712	91	1778556	[putative cobalamin synthesis protein (Escherichia coli)]	53	32	582
88	7	5217	4342	91	1208719	[putative fibrillar-associated protein (Heliobacterium salinarum)]	53	38	876
92	5	2395	2028	91	165366	[glutamate oxidoreductase (Gluconobacter oxydans)]	53	33	708
96	9	6832	7762	91	151204	[ORF1, putative 42 kDa protein (Streptococcus pyogenes)]	53	42	1131
100	8	7629	8600	91	149591	[maturation protein (Lactobacillus paracasei)]	53	32	972
128	9	6412	6972	91	17017237	unknown Mycobacterium tuberculosis	53	33	561
128	12	8429	9253	91	131070	[pentraxin fusion protein (Xenopus laevis)]	53	31	825
146	1	3	950	91	161607	[probable hemolysin precursor (Streptococcus agalactiae (strain 74-260))]	53	36	948
165	2	2162	3022	91	1755130	[inoturnin (Xenopus laevis)]	53	30	861
171	3	2304	2624	91	1732200	[PTS permease for maltose substrate (Hem IVBrio furmisi)]	53	32	321
182	5	3785	3051	91	1700572	unknown (Bacillus subtilis)	53	35	725
209	3	2948	1935	91	1778505	[feric acid transport protein (Escherichia coli)]	53	28	1014
218	5	3884	2406	91	170162	[murF gene product (Bacillus subtilis)]	53	34	1479
250	3	473	790	91	170134776	[YbaB protein (Bacillus subtilis)]	53	30	318
275	1	1	1611	91	1701314	[YgeW (Bacillus subtilis)]	53	35	1611
302	1	544	2	91	140938	[murD (Bacillus subtilis)]	53	31	543
2	2	2543	3445	91	1701323879	[hypothetical protein (Bacillus subtilis)]	53	33	903
3	12	22402	23176	91	1701323879	[murF gene product (Mycobacterium neoaurum)]	52	36	975
5	3	8034	2356	91	1701324915	[lactate permease (Streptococcus anginosus)]	52	32	5739
22	16	19943	20123	91	170132501	[orf 3 (Spirochaeta aurantia)]	52	35	252
22	31	23140	24466	91	170132562	[ecmB ORF (Bacillus subtilis)]	52	32	1227
27	6	2357	4657	91	170135573	[P20 (AA 1-178) (Bacillus licheniformis)]	52	35	597

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession (nt)	Match gene name	% sim	% ident	Length (nt)
35	10	8604	7357	gi 508241	putative O-antigen transporter [Escherichia coli]	52	27	1248
45	4	4801	3662	gn PID j0102243	[AB005554] homologs are found in E. coli and H. influenzae; see SWISS PROT AC007420 [Bacillus subtilis]	52	36	1140
48	18	14385	13726	gn PID j205174	orf2 [Lactobacillus helveticus]	52	35	664
49	4	3321	5755	gi 1317740	[AF013987] nitrogen regulatory IIA protein [Vibrio cholerae]	52	19	435
54	4	2773	4668	gi 1506472	M. jannaschii predicted coding region M1577 [Methanococcus jannaschii]	52	36	1896
54	6	3250	4969	gi 12184853	(AB000797) Y410 [Rhizobium sp. NGR234]	52	40	282
66	6	8400	6955	gi 143140	TAG protein [Escherichia coli]	52	30	1466
71	26	10693	11132	gn PID j014993	unknown [Mycobacterium tuberculosis]	52	23	654
75	2	1673	1035	gn PID j0102271	[AB001683] Fara [Streptococcus sp.]	52	27	639
81	3	1439	2893	gn PID j011458	translucase kinase [Bacillus subtilis]	52	32	1455
81	8	4987	5781	gi 147403	mannose permease subunit II-p-han [Escherichia coli]	52	37	795
83	21	20687	21853	gi 143165	phosphoribosyl aminoimidazole carboxylase II (Pur-R; tgg start codon) [Bacillus subtilis]	52	37	1167
86	6	5785	4592	gi 1276879	bpaf [Streptococcus thermophilus]	52	26	1194
86	20	11620	17055	gi 454444	orf3 [Schistosoma mansoni]	52	26	1530
96	13	10540	9659	gi 288299	orf1 gene product [Bacillus megaterium]	52	33	882
111	1	2	2026	gi 148309	cytolysin A transport protein [Enterococcus faecalis]	52	27	2025
112	2	1457	2167	gi 471234	orf1 [Haemophilus influenzae]	52	33	711
118	3	2931	2385	bss151233	orf1 [Haemophilus influenzae]	52	33	567
122	9	5646	5581	gi 8214	Mp24 kDa macrophage infectivity potentiator protein [Legionella pneumophila, Philadelphia-1, Peptide, 184 aa] [Legionella pneumophila]	52	33	567
122	11	6359	6374	gi 444025	myosin heavy chain [Drosophila melanogaster]	52	36	306
122	11	6359	6374	gi 444025	dihydrolypamide acetyltransferase [Polibacter carolinicus]	52	52	216
134	6	4880	6313	gi 153733	M protein trans-acting positive regulator [Streptococcus pyogenes]	52	43	1314
135	3	7258	2716	gn PID j040024	unknown [Mycobacterium tuberculosis]	52	35	1079
141	3	1481	2319	gn PID j0100573	unknown [Bacillus subtilis]	52	35	639
161	4	2562	3024	gi 146241	22.4% identity with Escherichia coli DNA-damage inducible protein ... [Escherichia coli]	52	36	2463
173	2	958	183	gi 1215693	putative orf1; orf2 [Mycoplasma pneumoniae]	52	30	786

TABLE 2  
S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match	Gene name	% sim	% ident	Length (nt)
198	6	4400	3567	gi 191031301	hypothetical protein [Bacillus subtilis]		52	26	834
210	12	8844	9107	gi 497647	RNA gyrase subunit B [Mycobacterium genitalium]		52	38	264
214	10	5264	5431	gi 550697	envelope protein [human immunodeficiency virus type 1]		52	36	168
225	1	15	884	gi 1552773	hypothetical [Escherichia coli]		52	34	870
230	1	35	342	gi 10100582	unknown [Bacillus subtilis]		52	31	284
287	1	871	2	gi 101335028	protease/peptidase [Mycobacterium leprae]		52	29	870
361	2	1305	4	gi 1913194	TM-291 membrane associated protein [Trypanosoma brucei]		52	32	1304
23	2	2048	1173	gi 1014254943	unknown [Mycobacterium tuberculosis]		51	30	876
29	3	742	1521	gi 1929900	5'-methylthioadenosine phosphorylase [Salmonella enteritidis]		51	31	780
45	1	410	1597	gi 1877429	integrase [Streptococcus pyogenes phage T12]		51	32	1188
48	26	19227	18946	gi 12314455	(A000033) transcrition [unc]		51	33	282
73	5	4276	4816	gi 1474177	alpha-D-1,4-glucosylase [Staphylococcus xylosus]		51	31	261
81	11	8935	12057	gi 1311070	putative fusion protein [Xenopus laevis]		51	31	3123
83	5	1195	1386	gi 101401316	Vqfi [Bacillus subtilis]		51	33	792
98	10	7531	8538	gi 141500	ppp 3' (AA 1-353) 3'8 10 [pct-1] [Escherichia coli]		51	28	1008
113	6	3908	5173	gi 1466882	ppa1: B1496_C2_189 [Mycobacterium leprae]		51	27	1266
124	1	326	57	gi 1231166	(A000270) contains similarity to myosin heavy chain [Arabidopsis thaliana]		51	32	270
139	10	7286	6816	gi 1046241	orf14 [Bacteriophage HP3]		51	30	471
143	3	4983	3993	gi 11254935	probable copper-transporting apase [Escherichia coli]		51	26	981
148	15	11359	10226	gi 12293256	(A008220) putative hippurate hydrolase [Bacillus subtilis]		51	36	1134
149	8	6003	7113	gi 1033572	Herpesvirus saimiri ORF7 homolog [Kaposi's sarcoma-associated herpes-like virus]		51	21	1311
151	9	12052	115550	gi 1014281580	hypothetical 40.7 kd protein [Bacillus subtilis]		51	34	543
159	6	7555	1208	gi 116947	non-acetylneuraminic acid synthase [Escherichia coli]		51	36	654
174	1	1797	4	gi 1773166	probable copper-transporting apase [Escherichia coli]		51	28	1794
245	22	323	1753	gi 1014281580	anti-P. falciparum antigenic polypeptide [Samir acturem]		51	18	459
277	2	443	1311	gi 1532915129	p110 protein - Malaria gonorrhoea		51	33	669

TABLE 2  
5. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	length (nt)
350	1	890	3	gi 290909	oj307 [Escherichia coli]	50	50	888
363	4	1228	4485	gi 1707247	partial CDS [Caenorhabditis elegans]	50	50	3258
367	1	1701	4	gi 393394	786-791 membrane associated protein [Pseudomonas brassicae subgroup]	51	32	1898
15	5	5174	4497	gi P101458151	P3 [Bacillus subtilis]	50	38	678
16	4	2220	2882	gi P101425010	hypothetical protein [Bacillus subtilis]	50	29	363
19	5	2591	4159	gi 1552723	similar to voltage-gated chloride channel protein [Escherichia coli]	50	30	1269
25	4	2701	1997	gi 887860	ORF 4219 [Bacillus subtilis]	50	27	705
35	1	211	417	gi P1014236897	unknown [Saccharomyces cerevisiae]	50	33	207
39	4	3416	5152	gi P101420974	unknown [Bacillus subtilis]	50	27	1737
51	7	4000	5181	gi 1592027	carbamoyl-phosphate synthase, pyrimidine-specific, large subunit [Methanococcus jannaschii]	50	27	1182
51	9	7179	8303	gi 1591847	type I restriction-modification enzyme, S subunit [Methanococcus jannaschii]	50	28	1125
52	8	8740	9531	gi 144397	acetyl esterase [yinc] [Caldococcus saccharolyticum]	50	34	795
52	16	16591	15770	gi 2108229	basic surface protein [Lactobacillus fermentum]	50	34	822
57	7	4031	630	gi 152525	60S ribosomal protein L7B [Schistosoma carnosus ponae]	50	40	306
71	23	29148	28983	gi P101401328	YqjA [Bacillus subtilis]	50	30	965
86	12	11155	10769	gi P101424964	hypothetical protein [Bacillus subtilis]	50	24	187
93	2	1205	330	gi 1066016	similar to Escherichia coli pyruvate, water dikinase, Swiss-Prot Accession Number P23538 [Pyrococcus furiosus]	50	24	876
96	5	1673	2959	gi P101422433	glutamate:pyruvate synthetase [Staxosira juncea]	50	29	1287
96	2	218	1171	gi 151110	[leucine-, isoleucine-, and valine-binding protein [Pseudomonas aeruginosa]	50	30	954
103	4	3303	2785	gi 1155250	O-antigen ligase [Salmonella typhimurium]	50	31	519
115	5	6480	5980	gi 895247	putative cbl operon regulator [Bacillus subtilis]	50	26	503
139	1	7529	7095	gi 1121675	skelatal muscle ryanodine receptor [Homo sapiens]	50	32	255
129	13	8192	7865	gi 152271	319-kDa protein [Rhizobium meliloti]	50	30	221
131	5	7834	6819	gi 140348	prot. resolvease Tmp I (AA 1 - 284) [Bacillus thuringiensis]	50	35	816
153	1	1	597	gi P101402015	similar to NITROGENASE, [Bacillus subtilis]	50	29	597

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start	Stop	match accession	search gene name	% sim	% ident	Accession (nt)
155	5	5986	5973	gi 13748693	Spac (Streptococcus thermophilus)	50	28	555
160	9	7190	6323	gi 1786983	(AE000179) g311: 92 pct identical to the 313 aa hypothetical protein YNHE_ECOLI SW: P51697; 26 pct identical (7 gaps) to 167 residues of the 373 aa protein nuc-TRICU SW: P48057; SW: P52697 (Escherichia coli)	50	30	1068
163	6	7386	8091	gnl P10101313	Ypsk (Bacillus subtilis)			
167	6	5232	3340	gi 1413526	lipo-2r gene product (Bacillus subtilis)	50	27	1293
169	2	1030	1340	gnl P10104540	endolysin (Bacteriophage Bacillus)	50	35	576
171	5	3168	4025	gi 1606080	ORF_0200, Genomeplot suggests frameshift linking to 0267, not found	50	27	858
210	11	8121	814	gi 130038	HRV 2 polypeptide (Human rhinovirus)	50	25	264
364	7	1538	135	gi 1351322	Pro-292 membrane associated protein (Trypanosoma brucei subgroups)	50	31	1404
10	7	5911	5090	gi 144659	ORF B (Clostridium parvifungens)	49	24	822
26	5	10754	9768	gi 142410	ATP-dependent nuclease (Bacillus subtilis)	49	31	987
66	7	9777	8398	gi 141470	Irka gene product (Methanosarcina mazei)	49	26	1380
77	6	5265	4648	gnl P1010283322	Peck protein (Mycobacterium avium)	49	28	717
82	13	12689	13249	gnl P101025591	Hypothetical protein (Bacillus subtilis)	49	20	561
91	9	4866	4531	gi 140067	X gene product (Bacillus sphaericus)	49	26	336
112	5	4019	4548	gi 1574380	Ilc-1 operon protein (IlcB) (Haemophilus influenzae)	49	27	730
129	7	6098	4949	gnl P1010267587	Unknown (Bacillus subtilis)	49	35	1110
135	5	3875	4318	gi 139573	P20 (AA 1-178) (Bacillus licheniformis)	49	25	584
154	2	1423	1953	gnl P101010102	Regulatory components of sensory transduction system [Synchocystis sp.]	49	29	531
156	5	3508	1637	gnl P10101732	Hypothetical protein (Synchocystis sp.)	49	25	1242
173	5	3500	2940	gi 1400324	LORF X gene product (unidentified)	49	30	561
182	1	1057	2	gi 131002	First methionine codon in the 80251 ORF (Saimirine herpesvirus 2)	49	25	1056
192	6	5352	3667	gi 1294472	(AF034499) contains similarity to homeobox domains (Caenorhabditis elegans)	49	23	1685
253	4	1129	1350	gi 151116	SR44 protein (Saccharomyces cerevisiae)	49	23	222
277	1	600	136	gi 136864	ORF (18 kDa) (Vibrio cholerae)	49	32	465
322	1	1435	887	gi 1733524	phosphatidylinositol 4,5-bisphosphate 3-kinase (Dictyostelium discoideum)	49	24	549



TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Match accession	Match gene name	% sim	% ident	length (nt)
365	3	1436	132	gi 391394	Tb-291 membrane associated protein (Trypanosoma brucei subgroup)	49	31	1005
33	7	4461	3277	gi 145844	codes for a protein of unknown function [Escherichia coli]	48	26	1185
40	2	652	1776	gnl P101430649	ornithine decarboxylase [Nicotiana glauca]	48	29	1125
67	4	1377	2384	gi 1772652	2-keeto-3-deoxygluconate kinase [Haloferax volcanii]	48	30	1008
74	2	4369	3871	gi 2183678	(AE000101) Y402 HNLobium sp. NC2341	48	27	399
81	2	1326	541	gi 1531672	[lactose repressor (Streptococcus mutans)]	48	33	786
81	4	2981	3646	gi 146042	[fructose-1-phosphate aldolase (E. coli) [Escherichia coli]]	48	30	666
97	1	602	51	gi 153794	[reg (Streptococcus gordonii)]	48	29	552
110	1	1	3132	gi 1181114	[p18 gene product (lactococcus lactis)]	48	23	3132
131	5	2914	2147	gnl P10183811	[acyl-ACP thioesterase [Brassica napus]]	48	27	768
133	4	1894	2628	gnl 1541988	[putative ORF (Bacillus subtilis)]	48	27	867
139	6	4231	4599	gi 1049388	[ZK470.3 gene product (Caenorhabditis elegans)]	48	23	369
139	8	5016	5645	gi 1027732	[unim (Staphylococcus haemolyticus)]	48	29	630
140	12	11936	11007	gnl P10102049	[H. influenzae, ribosomal protein almine acetyltransferase: 241105 (189) (Bacillus subtilis)]	48	27	910
146	9	5670	4654	gi 1591731	[polyomate kinase (Methanococcus jannaschii)]	48	24	1017
161	3	1280	2374	gnl P10101578	[collagenase precursor (EC 3.4.-.-) [Escherichia coli]]	48	24	1095
172	11	10581	11048	gnl P10101132	[hypothetical protein (Synecocystis sp.)]	48	27	468
182	4	2930	2586	gi 140067	[X gene product (Bacillus sphaericus)]	48	37	345
210	15	10796	11196	gnl 138401028	[LARGE OSMOTIC STRESS-INDUCED PROTEIN D-29 (LEX D-29)]	48	30	411
214	12	6231	6482	gi 140389	[non-toxic components (Clostridium botulinum)]	48	26	252
221	10	1000	1000	gi 1257364	[H. influenzae predicted coding region H0392 (Haemophilus influenzae)]	48	27	702
227	2	647	3928	gi 1673693	(AE000005) Mycoplasma pneumoniae, COG_0718 Protein (Mycoplasma pneumoniae)]	48	30	3382
253	2	480	758	gnl P101018697	[unknown (Escherichia coli)]	48	31	279
363	3	1874	1122	gi 181137	[cgr-4 product (Chlamydomonas reinhardtii)]	48	40	753
369	1	905	2	gi 181137	[cgr-4 product (Chlamydomonas reinhardtii)]	48	38	504
3	21	20879	22258	gnl P1010264778	[putative maltose-binding protein (Streptococcus coelicolor)]	47	33	1180

TABLE 2  
S. pneumoniae - putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
6	4	4089	4598	gi135513	120 (AA 1-119) [Bacillus licheniformis]	47	23	570
15	3	3736	1760	gn1210d100572	unknown [Bacillus subtilis]	47	25	1977
35	15	14516	13263	gi1173351	Cap5L [Staphylococcus aureus]	47	20	1254
51	6	3547	4002	pir1A3024/A370	12K antigen precursor - Mycobacterium tuberculosis	47	38	456
55	6	10134	9273	gi139648	US [Bacillus subtilis]	47	26	882
92	4	1753	3276	gn1210d100611	PCPC [Streptococcus pneumoniae]	47	35	1524
127	9	5589	5386	gi1786458	(AP000134) f120; This 120 aaorf is 76 pct identical to g9991 to 42 residues of an approx. 48 aa protein f127_UBIN Sm: P4395 [Bacillus coli]	47	32	204
130	2	1232	1759	gn1210d100655	unknown [Mycobacterium tuberculosis]	47	23	528
140	4	4951	3542	gn1210d100964	homologue of hypothetical protein in a tetracycline synthesis gene cluster of Streptomyces hygroscopicus [Bacillus subtilis]	47	24	1410
151	4	4814	4200	gi1322674	M. jannaschii predicted coding region MJC141 [Methanococcus jannaschii]	47	27	615
157	3	803	1174	gn1210d101320	Vag2 [Bacillus subtilis]	47	25	372
178	5	3267	2155	gi1246730	(AP000307) 1337; sequence change joins ORF 938 & 935 from earlier version (V03K_ECOLI Sm: P42539 and V03K_ECOLI Sm: P42609) [Bacillus coli]	47	30	1113
273	1	2	1549	gn1210d1005973	autocatalytic sensor kinase [Bacillus subtilis]	47	32	1548
300	2	880	644	gi13035755	zinc finger protein Png-1 [Mus musculus]	47	22	237
54	14	14182	12638	pir1A345601/S436	rota protein - Streptococcus pyogenes	46	24	1545
88	1	2	1018	gn1210d1023891	xylose repressor [Maerococcus thermophilus]	46	27	1017
96	7	4533	5860	gn1210d101652	ORF_ID:034715; similar to [UniseqPro Accession Number P45721] [Bacillus coli]	46	23	1308
112	1	1127	3	gi12209215	(AP004325) putative oligosaccharide repeat unit transporter [Streptococcus pneumoniae]	46	24	1125
122	13	7098	7982	gi1054776	bica4 gene product [Bacillus subtilis]	46	34	675
127	14	2198	8125	gi1469286	afuA gene product [Actinobacillus pleuropneumoniae]	46	28	1074
132	4	7033	6197	gi1153794	508 (8a) [Bacillus subtilis]	46	26	897
140	8	4220	7723	gi1235795	pullulanase [Thermomonas thermophilus]	46	21	498
140	9	9205	8315	gi1407878	leucine rich protein [Streptococcus agalactiae]	46	27	891

TABLE 2

5. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	Accession	Match gene name	% sim	% ident	length (nt)
162	1	1	1125	gi 1143209	ORF7, Method: conceptual translation applied by author [Shigella sonnei]	46	25	1125
199	1	1	585	gi 1547171	(AF000296) No definition line found (Candida albicans)	46	28	585
233	3	1971	1477	sp P02562 MYSS	MYOSIN HEAVY CHAIN, SKELETAL MUSCLE (FRAGMENTS)	46	27	495
232	2	760	1608	gi 1016112	MYF18 gene product (Cyprinodon paraxodon)	46	28	849
292	1	987	220	gi 1673744	(AF000011) Mycoplasma pneumoniae, cytidine deaminase; similar to GenBank Accession Number C5312, from M. pirum [Mycoplasma pneumoniae]	45	29	458
30	8	5843	6472	gi 1768049	(AF000270) o235; This 215 aa orf is 79 pct identical (10 gaps) to 196 residues of an approx. 216 aa protein YTB_BAC20 SM: P05568 [Escherichia coli]	45	24	630
48	6	3461	3868	gi 722339	unknown [Acetobacter xylinum]	45	29	408
60	1	307	2	gi 1699079	coded for by C. elegans cDNA yk41h4.3; coded for by C. elegans cDNA yk15d49.5; coded for by C. elegans cDNA yk59a10.5; coded for by C. elegans cDNA yk41h1.3; coded for by C. elegans cDNA cn20g10; coded	45	36	306
72	16	1437	1484	gi 132100	MDH (dehydrogenase (ubiquinone) [Artemia franciscana])	45	25	504
99	7	9158	7941	gi 152192	mutation causes a succinylglutamate-phenotype; ExoQ is a transmembrane protein; third gene of the exoQ operon; putative [Rhizobium meliloti]	45	28	1218
127	12	7046	6606	sha 153689	HAIR-iron utilization protein [Haemophilus influenzae, type b, H42, HNT 79106, Peptide, 506 aa] [Haemophilus influenzae]	45	24	441
137	5	1561	2619	gi 472921	gamma-type Na-ATPase [Enterococcus faecalis]	45	33	1059
209	1	774	364	gi 104141	restriction endonuclease beta subunit [Bacillus coagulans]	45	20	411
314	1	604	2	gi 1480457	latex allergen [Hevea brasiliensis]	45	31	603
20	18	19782	20288	gi 433942	ORF [Lactococcus lactis]	44	26	507
87	8	7030	6452	gi 537207	ORF_E277 [Escherichia coli]	44	26	579
166	5	4909	4037	gnl P1D130802	membrane transport protein [Bacillus subtilis]	44	25	873
247	1	818	75	gnl P1D100716	ORF1 [Bacillus sp.]	44	20	744
32	3	1885	3876	gi 2151768	Pyph [Streptococcus pneumoniae]	43	24	1992
36	17	15467	18276	gi 1045739	M. genitalium Predicted coding region M2064 [Mycoplasma genitalium]	43	26	2790
54	15	14556	17343	gi 520941	penicillin-binding protein 1A and 1B [Bacillus subtilis]	43	27	2688
67	2	696	1352	gi 536934	y3cA gene product [Escherichia coli]	43	29	657
139	2	2416	338	gi 1394600	similar to eukaryotic Na/H exchangers [Escherichia coli]	43	24	2079

TABLE 2

S. pneumoniae - Putative coding regions of novel proteins similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)	match accession	match gene name	% sim	% ident	length (nt)
298	1	3	809	gi 411972	lps-48, gene product [Bacillus subtilis]	43	24	807
387	1	47	427	gi 2115652	[AP016669] No definition line found [Caenorhabditis elegans]	43	30	381
385	4	3231	3127	gi 2142399	[AF000073] Y4EP Phalloidin sp. R0234	41	25	1095
340	1	582	70	gn PTD e218681	CDP-diacylglycerol synthetase [Arabidopsis thaliana]	41	20	513
386	1	805	914	gi 1156742	P27-2 protein [Trypanosoma cruzi]	41	27	2292
368	2	2	843	gi 211783	LMW glutenin (AA 1-356) [Triticum aestivum]	41	34	942
335	3	4469	2861	gi 42023	Members of ATP-dependent transport family, very similar to mdr proteins and multidrug resistance proteins [Homo sapiens]	40	18	1629
365	2	55	1438	gi 1033972	Herpesvirus saimiri ORF73 homolog [Kaposi's sarcoma-associated herpes-like virus]	40	21	1344
1	3	2979	3860	gn PTD d101908	hypothetical protein [Synechocystis sp.]	39	26	882
3	5	3814	4847	gn PTD d101961	hypothetical protein [Synechocystis sp.]	39	19	814
26	6	14035	10724	gi 142439	ATP-dependent nuclease [Bacillus subtilis]	38	20	1212
47	1	3	4916	gi 612549	NP-180 [Pentatomys marinus]	36	23	4914

TABLE 3

S. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
1	4	1428	3009
1	6	4611	4964
3	2	818	994
3	3	1182	1574
3	7	5382	6497
3	25	25046	25396
3	26	25625	26317
6	2	1519	1689
6	14	12875	12618
6	15	13215	12841
6	18	12577	15390
7	12	9955	9419
7	13	10361	9910
8	6	3915	4280
9	9	6024	5704
10	8	6900	6288
10	9	7136	6888
11	1568	7672	7248
12	1	1140	4
12	3	1779	1456
14	2	1913	1434
16	1	1	243
16	5	5675	3087
17	1	324	34
17	3	1451	1050
17	9	4890	4465
20	14	14544	15893

TABLE 3

S. pneumoniae Putative coding regions of novel proteins not similar to known proteins

Config ID	ORF ID	Start (nt)	Stop (nt)
21	3	3359	2589
21	5	4802	4482
22	21	17099	17462
22	25	19467	19982
22	33	25540	25744
22	35	26388	26218
22	36	26382	27572
23	7	6655	6032
23	8	7112	6653
24	1	36	518
25	3	3005	2641
27	4	4819	4223
27	5	4789	4956
28	5	3017	1797
28	8	4272	3850
28	10	5028	4597
28	11	5746	5072
29	7	5596	4919
29	8	5039	5518
29	9	5595	8207
30	9	6511	6263
31	6	2664	2144
32	5	5203	5538
33	8	5127	4668
34	10	8024	7740
34	12	9360	8643
34	13	9667	9377

TABLE 3

5. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
34	18	113104	11902
35	11	9648	8588
35	12	11073	9670
36	2	334	1041
36	12	11120	10853
36	13	10993	11388
36	15	12172	14595
38	7	4269	4577
38	8	4480	5001
38	10	5517	5711
38	17	10732	11376
40	3	1728	3143
43	1	172	5
43	7	8884	8732
43	8	9568	9071
44	4	6831	6831
45	3	3204	3665
46	4	3875	3468
46	7	6074	7081
48	5	3196	3582
48	8	4579	4229
48	11	9323	8922
48	16	13042	12494
48	20	16342	15764
48	24	17971	18351
48	30	21979	21776
49	1	209	3

TABLE 3

S. pneumoniae - Punctate coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
50	4	3307	2672
51	5	3239	3598
52	11	12146	11883
54	7	5888	5187
54	8	6011	5939
54	9	6004	6210
54	16	17688	18068
55	9	10515	10123
55	12	11947	11241
56	3	935	1387
56	4	1496	1939
57	3	1624	2130
57	4	2100	2501
58	6	7541	7335
59	1	2	430
59	4	2416	2736
59	5	2734	3063
59	8	4743	5549
59	9	5459	5929
60	6	5741	6451
61	3	2395	1772
61	5	3316	3176
64	1	2722	2
66	2	1180	3147
66	8	9082	9495
67	3	1143	1182
69	2	1165	980



TABLE 3

S. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
70	5	4059	3922
70	6	4215	4057
70	9	5268	5504
71	15	20351	21901
71	16	21859	22338
71	19	26204	27556
72	9	8458	6081
73	6	3815	4216
73	7	4369	4773
73	10	7183	6428
73	15	9462	9668
76	1	524	195
76	2	867	535
76	11	8602	9210
80	6	7924	8109
81	1	244	2
81	10	6631	8931
81	14	1072	1150
83	17	16610	16460
84	3	4684	2929
86	2	2147	1092
86	4	3606	2875
86	19	16767	17114
87	5	5326	5000
87	7	6459	6003
87	9	7224	7006

**TABLE 3**  
 S. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
	87	18	17530
	87	19	18275
	88	2	1619
	88	4	2711
	88	9	6252
	89	9	7371
	90	2	899
	90	3	1143
	91	3	2859
	91	4	3170
	91	6	4253
	93	1	391
	93	6	2668
	93	8	4533
	96	1	3
	96	2	904
	96	3	1407
	96	4	1250
	97	9	7043
	99	15	18522
	99	17	39717
	100	2	4094
	103	1	48
	103	6	4924
	104	5	6142
	105	7	6098

TABLE 3  
 S. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
106	1	1	363
106	10	9832	10212
108	1	2	268
111	3	3417	3788
111	4	3809	4606
115	10	10854	110438
116	3	2873	2121
118	2	2274	1357
122	4	2598	2333
122	10	5858	6199
122	12	6301	7116
124	2	346	690
128	4	2444	3368
129	1	689	102
129	2	1011	178
129	8	6454	6056
129	9	6340	6277
129	12	7809	7821
131	3	1435	756
131	10	5972	5673
134	11	11838	11209
135	2	625	1140
136	4	2913	3830
137	2	325	134
139	12	14027	14521
139	13	14840	14932
139	14	15363	14875

TABLE 3

S. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
140	20	19822	20038
142	1	1	285
146	3	760	479
146	4	1149	778
146	7	3604	2885
146	13	8223	8403
146	14	9399	10676
146	15	110952	9750
147	7	7488	7276
147	9	8913	866
148	7	5298	4765
149	1	2	1336
149	3	2557	2880
149	9	6586	6070
150	2	1355	579
150	3	2556	1909
153	3	2061	2642
154	3	1953	1741
155	2	2181	1411
156	8	4550	4311
157	1	37	284
159	2	631	780
159	4	1384	1722
159	7	3271	4017
161	2	1332	1018
165	3	5535	6945
166	6	5406	6972

TABLE 3

S. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (aa)	Stop (aa)
167	9	6075	6395
169	5	2828	3205
170	7	6485	6243
170	8	5964	5362
170	9	7303	5962
170	11	8790	7806
171	9	7150	7476
172	5	2398	1848
173	4	2913	2677
175	2	639	835
175	3	893	1789
176	2	1487	546
176	3	2200	1466
177	9	4886	4925
177	10	4923	5177
177	11	5111	5347
177	13	7396	6703
178	6	3452	3724
181	5	1853	2473
182	2	2112	1102
182	3	2617	2006
183	2	2126	2320
185	5	6893	4219
185	6	4846	4634
187	4	2540	3557
188	4	3686	4363
188	5	4181	4821

TABLE 3

5. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
188	6	5882	6493
189	5	3143	2844
189	9	5956	5544
191	1	618	4
191	11	1035	1000
192	3	2861	2268
192	4	308	
192	7	6800	5311
193	3	997	835
194	4	2315	2127
195	5	6249	4583
195	6	6620	6231
196	2	1553	1849
197	1	1	861
198	9	6844	6644
200	5	5329	5769
200	6	5993	6595
204	5	3914	3276
205	2	447	1709
209	4	2038	2460
209	5	2458	2682
210	10	7376	6230
210	13	9029	10441
210	14	1035	10705
214	5	2581	2330
214	9	5085	5277
214	11	5996	5754

TABLE 3

5. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
217	2	541	194
218	2	914	1432
218	3	1430	1972
218	6	3639	3821
219	1	458	39
220	1	669	600
223	4	2617	1964
227	1	1	510
234	4	1539	1312
234	6	2116	1838
235	1	52	312
235	2	310	63
238	1	660	64
246	1	1	20
248	1	3	362
248	2	443	1222
254	3	2789	792
258	2	1179	1616
260	3	1770	2123
263	1	653	177
263	4	2244	1900
263	5	3969	2973
266	1	1	342
266	2	177	1022
270	2	1124	1681
272	1	857	186
275	2	1484	2295

TABLE 3

S. pneumoniae - Putative coding regions of novel proteins not 4100<sub>aa</sub> to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
278	1	2	405
282	1	714	391
282	4	1463	1114
287	2	1119	826
288	1	500	4
289	1	684	4
293	5	1585	1858
293	2	2539	2925
294	1	21	608
296	2	494	700
296	3	670	843
302	1	261	530
309	3	559	350
310	2	249	1889
316	2	2087	1818
317	2	1048	584
318	2	313	777
319	3	477	333
327	2	912	607
331	1	1	549
333	1	2	335
333	2	405	82
333	3	127	342
341	1	1	705
345	2	895	701
346	2	750	139
349	1	1	198



TABLE 3

S. pneumoniae - Putative coding regions of novel proteins not similar to known proteins

Contig ID	ORF ID	Start (nt)	Stop (nt)
350	2	81	413
355	1	44	973
358	2	636	446
360	2	548	628
362	2	1635	1265
378	1	345	1064
379	2	683	510
381	1	109	693
385	1	150	4
385	2	269	30

148

## (i) GENERAL INFORMATION:

(i) APPLICANT: Charles Kunsch  
Gil H. Choi  
Patrick S. Dillon  
Craig A. Rosen  
Steven C. Barash  
Michael R. Fannon  
Brian A. Dougherty

(ii) TITLE OF INVENTION: Streptococcus pneumoniae Polynucleotides and Sequences

(iii) NUMBER OF SEQUENCES: 391

## (iv) CORRESPONDENCE ADDRESS:

(A) ADDRESSEE: Human Genome Sciences, Inc.  
(B) STREET: 9410 Key West Avenue  
(C) CITY: Rockville  
(D) STATE: Maryland  
(E) COUNTRY: USA  
(F) ZIP: 20850

## (v) COMPUTER READABLE FORM:

(A) MEDIUM TYPE: Diskette, 3.50 inch, 1.4Mb storage  
(B) COMPUTER: HP Vectra 486/33  
(C) OPERATING SYSTEM: MSDOS version 6.2  
(D) SOFTWARE: ASCII Text

## (vi) CURRENT APPLICATION DATA:

149

(A) APPLICATION NUMBER:

(B) FILING DATE:

(C) CLASSIFICATION:

(vii) PRIOR APPLICATION DATA:

(A) APPLICATION NUMBER:

(B) FILING DATE:

(viii) ATTORNEY/AGENT INFORMATION:

(A) NAME: Brookes, A. Anders

(B) REGISTRATION NUMBER: 36,373

(C) REFERENCE/DOCKET NUMBER: PB340P1

(vi) TELECOMMUNICATION INFORMATION:

(A) TELEPHONE: (301) 309-8504

(B) TELEFAX: (301) 309-8512

150

(2) INFORMATION FOR SEQ ID NO: 1:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 5625 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 1:

CCAGCAAAA	CCAGCTACAG	CTAAAGGAAC	TTACGTAACA	AACTTGACTA	TCACAACTAC	60
TCAGGGTGT	GGTATCAAAG	TTGACGTAAA	CTCAGTTTAA	TCAGTAGTTA	AAGTAATGTA	120
AAAAAGTTGA	AGACGCTATG	TCTCAACTTT	TTTTGATGTA	CGACGGCAT	GTGTATAGT	180
AGATGTGTAC	TATTTCTAGT	TCAATCTACT	ATAGTAGCTC	AGAAGTCGGT	ACTTAAACGT	240
GCTATATCAA	AACCACTCCT	TGAAAAACGT	GGACTGGTTT	CGTGTTTGGA	TTATTACCTT	300
GAACGACATG	CGTTAAAAAT	TAGTTGAACC	GCCGTATGCC	GAACGACGT	ACGGTGGCT	360
GAGAGGGGCT	AGAGATTATC	CCTACTCGA	TTTCGAAATC	TAGTGAATG	AATCTGGAAT	420
AGTCCATCGA	GCTTCTAAT	ACTCTTCGAA	AATCTCTTCA	AACCACTCA	ACGTCGCCCT	480
GCCGTGCGTA	TGGTTACTGA	CTTCGTCACT	TCATCCACA	ACCTCAAAAC	AGTGTTTTGA	540
GCTGACTACG	TCAGTTCCAT	CTACAACTC	AAAACAGTGT	TTTGAGCAAC	CTGCGGCTAG	600
TTTCTAGTT	TGCTCTTTGG	TTTTCAATTGA	GTATAACACA	TTGTTAGAAG	TTGGTTTAAA	660
TTTCCTAATC	AGTTTGTCCA	CATTACCTT	CGATATATTA	TATCCCATAG	TTAAGGTGG	720
TCATACAGAT	GATTATAGTC	ATGAGCCGT	AAAACCTAGT	GTTCCTTTAG	TTGACAAAGA	780
TGCCATGAAA	AAAATATTTC	TAACTGTAA	AGGATATTTT	GAAATAAATA	TAGATGAJAA	840
TATCACCGAT	ATCTATACG	TAAATGGTAC	TGCTATTCCT	TATCTTTATT	TACGTTCAAT	900
TGTTTCAATA	GTTCGGCAA	TTGATAGCAG	TGAAGCAATG	TTGTACCTTA	TCATTAAATGT	960
TTTAGAGTTA	CTAGATAAAT	CTCAACCTTT	TGAAGAGAA	TAATTTATTA	GCTCACATAA	1020
TTGAGGGTAA	GGAAAAGTAA	AAGCAGTAAG	AAAAATGTCT	TGCATTATAC	AGCAACCTTT	1080
TGGGAATGAG	TGGATGGATT	GAATAAAATT	TGATTAAGAG	TGGATGATT	ATCTGTAGAT	1140
TATTATTGGA	CAGTTAGTCT	TGAAGTAGTC	TAGAATTAG	GTTATTAATCA	GTAGAAGCCT	1200
TGCTAATAAT	GAGGAGGTTA	GTTTATGTAT	AGTAGACTGA	ATCTAAATAA	GTACGAAACA	1260
ATTGCTAAAA	CATTATAGA	AATTAATTTT	ACTTCCCAA	TCGATTGTGT	CTCATCTTAT	1320
TTCAATCCGC	TATATATTA	GGTATCGAAT	CTTCATCAGA	ATGATAAAAT	TAATCAATTG	1380
ATATCTGATT	ACAAACAGAA	TATGAAAGCT	TTTTATATCA	CTATTGAAAA	ATTTATACGA	1440

GATGATGAA	GCCTTAACTG	TTATTTTATA	AAGGTTATTT	CAAGTCGTT	CAAGGTAACA	1500
AGTCTAGATC	AGATTGAAGC	TGATAAAACG	ATACAAAGAA	AATATTCAAG	TGAGCTAAAA	1560
AAATTTTATG	GATTTTATAA	TGAGATTATT	TGTGAAGAAA	ATAGTTTCCT	ACATGACGA	1620
AAGAGGTGCT	CGATTGCTT	TAGGTAGTCG	ATCCTGAGT	TGATAATTCT	CAAGGTATGG	1680
ACTTCTTTT	CATGAATCAG	GTAAAAGAGC	AGGTATTGTT	TAGAGACAAT	CATTCTGAGC	1740
ATATTTTCTG	GATAGAGGGA	GTATCCGATT	TTATGATCAA	AGTTAATACC	GCCTCTGGT	1800
GAGAAGATGA	GTAGGTGGT	AATTTAACT	ATTAAACAGA	ATTTTGTATT	AAAAGTATTA	1860
TTTCATGAGA	GAATCTTAA	TTTCACAATC	CATAGGCAAA	CGCTTGCAIT	TGTTTTTTTA	1920
TTGACTATA	ATAGGTTGCT	ATAAAGCCTT	CTGTATTAAT	AAAATGTAGA	AGGTGTAGAA	1980
AGTAAGGATT	TAGAAATATT	GTAGTTAAAA	ACACAATGTT	GCTATTCCIT	ACGATAGGGA	2040
GATAGATATG	CAATGATAG	AAGTGGAA	TCCTCAGAAA	AATTTTGTGA	AGACTGTTAA	2100
GGAACCGGG	TTGAAGGGG	CTTTCGCTC	CTTTATTTCAT	CCTGAAAGC	AGACCTTTGA	2160
AGCGGTCAAG	GATTTGACCT	TTGAGGTTCC	AAAAGGCGAG	ATTTTAGGAT	TTATCGGGGC	2220
AAATGGTGCT	GGGAATCGA	CAACCATTA	ATGCTGACA	GGAAATTTGA	AACCAACATC	2280
TGGTTTTGCT	CGGATTAAAG	GCAAGATTCC	CCAGGACAAT	CGGCAAGATT	ATGTCAAAGA	2340
TATTTGGCTA	GTCTTTGGAG	AACGCACCCA	GCTATGTTGG	GATTTGGCTC	TGCAGAGAC	2400
CTACACTGTC	TTAAAAGAGA	TTTATGATGT	GCCAGACTCG	CTCTTTCATA	AGCGTATGGA	2460
CTTTTGTGAT	GAGTCTTGG	AATTTGAAGGA	CTTTATCAAG	GATCCCTGTC	GGACTCTTTC	2520
ACTGGGACAA	CGGATCGGG	CGGATATTGC	GGCTCTCTTG	CTCCACAATC	CCAAGGTTCT	2580
TTTTTTAGAT	GAGCGGACCA	TTGGTTTGA	CGTTTCGGTT	AAGGATAATA	TTGTCGGGC	2640
AATTACTCAG	ATCAATCAAG	AGGAAGAAAC	TACCAATCTT	TTGACCACCT	ACGATTTTGA	2700
TGATATTGAG	CAACTTGTG	ATCGGATTTT	CATGATTGAC	AAGGGGCAAG	AGATTTTGA	2760
TGGAGCGGTG	AGCCAATCA	AGGAGACCTT	TGGTAAGATG	AAGACTCTCT	CTTTTGAAT	2820
GCTACCAAGT	CAAAATCATC	TGCTCTCTCA	CTATGACGGT	CTGTCTGATA	TGACCATTTGA	2880
TAGACAAGGA	AACAGCTTCA	ACATTGAATT	TGATAGTTCT	CGCTACCACT	CAGCTGACAT	2940
TATCAAGCAA	ACCTCTCTG	ATTTTGAAT	CCGCGATTG	AAGATGCTGG	ATACGATAT	3000
TGAGGATATT	ATCCGTGCT	TCTACCGAA	GGAGCTCTAG	GATGATCAAA	TTGTGGAGAC	3060
GTTATAAACC	CTTTTCAAT	GCAGGGGTTT	AGGAGTTGAT	TACTTACCGA	GTCAACTTTA	3120
TTCCTATATC	GATTTGGCAT	GTCAATGGGG	CTTTTGTGCC	CTTTTATCTC	TGGAAGGCTG	3180

152

TCCTTGATTC TCCGAAGAG TCCTTGATTC AGGGCTTCAG TATGGCGGAT ATCACCTCT	3240
ACATCATCAT GAGTTTGTG ACCAATCTTC TGACTAGATC CGATTCGTCC TTTATGATG	3300
GGGAGGAGGT CAAGGATGGC TCCATTATCA TGGCTTGTG GCGACCAATG CATTTGCGG	3360
CCCTCATCTC TTTCACCGAG CTGGTTCCA AGTGGTTGAT TTTTATCAGC GTTGCCTTC	3420
CATTTTAAAG TGTCATTTGC TTGATGAAAA TCATATCGG TCAAGGTATT GTAGAGGTGC	3480
TAGGATTAAC TGTCATTTAT CTTTTAGCT TAACGCTGCG CTATCTGATT AACTTTTCT	3540
TTAATATTG CTTTGGATT TCAAGCTTTG TGTTAAAAA TCTTTGGGT TCCAACCTAC	3600
TTAAGACTTC CATAGTGGCT TTTATGTCGG GGAGTTTGT TCCCTTGGCA TTTTTCCAA	3660
AGGTTGTTTC AGATATTCTC TCCTTTTGC CTTTTCATC CTGTATTAT ACTCCAGTTA	3720
TGATCATGT TGGAAAAATC GATGCCAGTC AGATTCTTCA GGCACCTCT TTGCAGTTCT	3780
TCTGGCTCTT AGTGATGGTG GGATTGTCTC AGTTAATTG GAAACGGCTC CAGTCCCTTA	3840
TCACCATCA AGGAGGTTAG TATGAAAAA TATCAACGAA TGCACTGTAT TTTTATCAGA	3900
CANTACATCA AACAAATCAT GGAATATAAG GTAGATTTTG TGCTTGTGT CTTGGGAGTC	3960
TTTCTGACTC AAGGCTTGAA TCCTTGTIT CTCAATGTCA TCTTCAACA TATTCATTTC	4020
CTAGAAGGCT GGACCTTTCA AGAGTAGCT TTCAATTATG GATTTTCCTT GATTCCTAAG	4080
GGAATGGACC ATCTCTTTT TGACAACTTC TGGGACTAG GCGAACGCT AGTCCGAAAA	4140
GGGAGTTTG ACAAGTATCT GACTCGTCCC ATCAATCCTC TCTTTCACAT CCTAGTTGAA	4200
ACCTTTCAGA TTGATGCTTT GGGTAACTC TTAGTCGGTG GTATTTTATT GCGAACACA	4260
GTGACCAACA TTGTTTGAC TCTTCCAAAA TTCTGCTTT TCCTAGTTTG TATTCCTTTT	4320
GCGACCTTGA TTTATACTTC TCTTAAAAATC GCAACAGCCA GTATCGCTT TTGACTAAG	4380
CAGTCAGGCG CCATGATTTA CATCTCTAT ACGTTCATG ACTTTGCTAA GTATCGATT	4440
TCATTTTACA ATTCTCTCT TCGTTGGTTG ATTAGCTTTA TCGTGCCTTT CGCCTTTACA	4500
GCCTACTATC CAGCTAGCTA TTTCTTACAG GAAAGGATG TGTCTTTAA CGTAGGAGGT	4560
TTGATGTGA TTTCTCTGTT TTTCTTTGTT ATTTCCCTTA AACTTTGGA TAAGGCTTA	4620
GATTCCTACG AAGTGCGGG TTCGTAAAAAG CTAAGTAAG ACTAAAAATCA AGAAAGAAC	4680
TTATGATGTT TGTAAATTGA GAAGTCAAG ATGAAAAATC AAAAAAGGCA GTTGTGCTG	4740
AGGTTTTGAA GGATTTGCCA GAATGTTTG GAATCCCAAG AAGCACACAA GCCTATATAG	4800
AAGGAACCA GACACTGCAA GTTTGGACCG CCTATCAGGA GAGTATTTG ACTAGATTG	4860
TAAGCTTATC CTATTCGAGT GAAGATTGTC CAGAGATTGA TTGTCTCGGC GTAAAAAGC	4920
TTATCAAGGT AGAAAAATTG GGAGCCAATT GCTTGCTACT TTAGAGAGTG AAGCTCGTAA	4980

153

AAAAGTGTG	TATCTGCAG	TCAAAACAGT	GCCAGAAGGT	TCTAATAAAG	ATTATGATCG	5040
AACAAATGAC	TTTTATCGAG	GTCTTGGCTT	TAAAAAGTTA	GAGATTTTTC	CTCAACTATG	5100
GAATCCGCAA	AATCCTGTGC	AGATTTTGAT	TAAAAAGCTT	GAATAATATT	ACTTGACATC	5160
TATTCTCAGA	GTGCTATACT	GTAAGTGTA	TCGCCGATTT	AGCTTAGTTG	GTAGAGCAAG	5220
GCACTCGTAA	AGCCTAGGTT	ATAGGTAGAT	AAAGCACTGA	GGATTTGAAA	AAATAGATAG	5280
GTAGAAGATA	ACCGTTAAGC	CTTACTCTTA	GCGGTTATTT	ATATTGTTTA	ATAGCGCTAA	5340
TATTTTATCA	ATTATGCGCTG	TTTTCGTGTT	TCITGGTAGTT	GTTCAGTTT	ATTGCTACTA	5400
TTTTTGATGG	TATGAATGTG	CTTATAATGT	ATCCCGGTTA	ACGAAAGTTT	TGGACTTATA	5460
CTCTTCGAAA	ATCTCTTCAA	ACCACCTCAA	CGTCGCCCTG	CCGTGCGTAT	GTTTATGACT	5520
TCGTCACTTC	TATCCACAAC	CTCAAAACAG	TGTTTTGAGT	GACTACGTCA	GTTCCTACTA	5580
CAACCTCAAA	ACACTGTTTT	GCCCAATCTG	CGGCTAGTTT	CCTAG		5625

(2) INFORMATION FOR SEQ ID NO: 2:

- (1) SEQUENCE CHARACTERISTICS:  
 (A) LENGTH: 7571 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 2:

CTCTCCAGCT	TTCCCTGCGA	GTTCGCCATG	TTGTGTCTTT	AAGAAGTCTA	AAAAATATCTC	60
CAATAAAACG	CATCGCTCTC	TCCTATCTCG	TTTCTCTGTG	TGAGTGTAC	TTGCCACAAT	120
GCTTACAAAA	TTATTACTT	TCTAGTCTGT	TAGGCTTGAG	GTTTCCGCTG	ATCTTGATTG	180
AATAGTTTCT	CGAACCACAA	ACCGCACAAAG	CTAGGCTTGC	TTTTTTTAGT	GCCATAACGC	240
CTCATCTTA	TCCATTATAA	CAAGAAAGCT	AGGCTTGAC	AAGCATCTTA	CGGAAATAGA	300
TTGACTATCG	AATCCCATAT	TGTTTGAGCC	TTTTCTTAA	TCCTCGCATC	TGAGATAGCC	360
CGGCTAGCCT	CATCTACTAG	ACTTTGCGCA	CGCCCTCGAA	TATCAGACAA	ATTATCATCT	420
GTCTGGCTAT	TATCATTGGT	TGTACTTGT	CTTTTGTAT	TGGCTGGTGC	AATTCATTT	480
TGCTTATAAG	CATTTTCAAC	CGTAAAGGTA	CTTCTGGCG	TATAAGGTAA	AATGGTATTG	540
GCAATGTTTC	TAAAGACHTG	AGCTGCACCG	TTTGAGTAG	AGCCAGCTAG	ATAGTGGTTT	600
TCATCACTGG	TCGAAAGGCC	AAGCCAGTGG	CTAATCACTA	CATCCGGAGT	ATAACCAATT	660
ACCCACTGGT	CACTTGTGTA	CTCCGGATTG	AAAACCTGCT	CAGTTGTTC	AGTTTCCCT	720

GCCATGACAT	AGTCTGCAGG	CGATGAACATA	ATACCGGTAC	CGTTGGTGAA	AGTCCCCAAC	780
ATCATACTGG	TCACTCTGTC	AGTACAGAC	TTATCAATCA	CCCCTTTTTC	TGA/TTTTTA	840
TGACTCGCAA	TAACTTGTC	ACTAGCATTT	TCAATTCTAC	TAATAAATG	AGCTTCAGGC	900
ATTAAACCTT	CATTTGCAAA	GGCGGCGTAT	GCTTGAGCCA	TTTGAAAGAG	GTTGGTTTCA	960
ACACCGCTTC	CCAAGGCGAC	ACCAAGAACA	CGGTGACCT	TTTCCATGTT	GAGTCCGAAT	1020
TTTTGCGCTG	CCTCAAAAGC	CTTGTCGACA	CCCAAAATCAT	TAACAGTGGC	AACAGCAGGT	1080
AGATTAAAGC	ATTCTGCCAA	GGCTTGATAC	ATAGGAACCT	CTCGACTCGT	TTTGATCCCT	1140
GCATAGTTAT	CAACCTTATA	GCTGTATATC	TGCATGGTAT	GGTTATCCAA	CTGCTTATTC	1200
AAAGCCGAGC	TTGCTTCAAC	TGCTGGCGTA	TAAACAACTA	AAGGCTTAAT	TGTAGAACCA	1260
GGACTAGCTT	TTGATTGGGT	TGCATAGTTG	AAATTCGGA	ATCCAGTTTT	ATCATTGTCA	1320
GCAACTTGAC	CGACAACCTC	ACGAACCTCC	CTGTGTTTCG	GTTTCGAGGC	TACACTTCCT	1380
GATTGAGCAA	ACGTTCCATC	CTCTGCCCTC	GGAAATAGCG	ATGTGTTTTT	ATAACAATTC	1440
TGTCATATTG	CTTGATAGTT	TTGGTCCAGC	TCTGTGTAAA	TGCGGTAGCC	ATTATTGACA	1500
ATCTCTTCCT	CTGTAGATT	ATACTTGGAA	ACAGCTTCAT	TAACCACGCG	ATCAAAATAA	1560
GAGGGGTAA	GGTAATCTGA	GATTTTTCCT	TCATACTTAT	CGTGCAATTG	CGAAGTCATA	1620
TCAACTTCAG	CAGCTTTGGT	TTCTTGGTTT	TTATCAATAT	ATCTCTGCTG	AACCATATTC	1680
TGCAAGACAG	TATCGCGCCG	ATTAGTAGAA	TCTTCTACGG	AACTCAAGGG	ATTATACAGT	1740
TCCGGCCCTC	TGAGCATCCC	TGCCAGATC	CGAGCTTGAT	CCAGACTCAC	TTCTGATGCA	1800
GAAACTCCAA	AGTATTTCTT	ACTCGCATCT	TCTACACCCC	ACACACCAAT	TCCAAAATAA	1860
CGCTGTGTAA	GGTACATGGT	TAGAATTTGC	TCTTACTAT	ATTTTTTTCCT	TAATTCTAAG	1920
GCAAGGAAAA	ATTCTTTCGC	TTTTCTCTCA	ACAGTTTGAT	CCTGCGATAA	ATAGGCGTTT	1980
TTAGCCAGCT	GTGCGTAAT	GGTAGAGCCA	CCACTGTGAC	GTCCAGCAGT	GACAAATAGC	2040
AAGAAAAAAC	GGCCATAGTT	AATCCCGTCA	TTTTTATAGA	AAGAACGGTC	TTCTGTGCGA	2100
ATAACAGCAT	TCTCAAGTT	TTTACTGTAT	TCAGTCAGCT	CAACATAGGT	TCCCTTTTGA	2160
CCAGACAAGG	CACCAGCCTC	TTTTCTTCTA	CGGTCAAAAA	TAAGAGTCCG	AGTTTTCAG	2220
GCATTTTGCA	AATCATTGAC	ATTGGTCGAC	TTGGCTACAG	CAACAAATA	GATTCCAACT	2280
AGCAAGCCTG	CACCTAAACC	TAGTATAAAG	ATAATCTTGT	TTAGATGATA	ACGACGCCAG	2340
AAATTTTGGAA	TCCGACCTAC	TTGGGCTAAT	TTTTTTCGAT	CACATACAGA	GCGACGTAAG	2400
ATAGTAGAAT	CAGAGTCTCT	TAGTTCACTT	GTTTCTTTTT	TAAAAAGAGA	AAGAAATTTT	2460
TCAATAAATT	TATCTAATTT	CATGCGTTTA	TTTTATCATC	TTTATCATAG	GAAGACAAGA	2520



ATTAGCTAT	TTCCTATCCA	AATAGGCTT	TTTTGTGTAC	AATATCTGTA	TGCAATTCAC	2580
ATTTACATTA	CCCSCCTCTC	TACCTCAAA	GACAGTAAAG	CAATTACTTG	AGGAACAACT	2640
CCTCATCCCT	AGAAAAATCC	GTCAATTTTT	GAGAATCAAG	AAACATATTT	TGATAAATCA	2700
AGAAGAGTC	CACCTGGANG	AAATCGTAA	TCTTGGAGAT	GTTTGCCAGT	TGACTTTTGA	2760
CGAGGAGAT	TATTCCTAAA	AGACGATCCC	TTGGGGCAAC	CCAGACTTAG	TGCAGGAAT	2820
TTATCAAGAT	CAACACTTGA	TTATTGTAAA	CAAAACCAGAG	GGGATGAAAA	CGCATGGTAA	2880
TCAACCAAA	GAAATTGCCC	TTCTTAAACA	TGTCAGTACC	TATGTGGGCC	AAACCTGCTA	2940
TGTCGTTCAT	CGTCGGACA	TGGAACCAG	TGGCTTAGTT	CTCTTTCCTA	AAAATCCTTT	3000
TATCCTCCCC	ATTCTCAATC	GCTTATTGGA	GAAGAAAGAG	ATTCTAGAG	AATATTGGGC	3060
TCTAGTTGAT	GGAAATATCA	ACAGAAAAGA	ACTTGTTTTC	AGAGACAJAA	TTGGACGTGA	3120
TGCGCATGAT	CGTAGAAAA	GAATAGTTGA	TGCAAAAAAT	GGGCAATATG	CTGAACGCA	3180
TGTAAAGCAG	TTAAAGCAAT	TCTCAAAACA	GACTTCCTTG	GCTCATTGCA	AGCTAAAGAC	3240
AGGGCGAACC	CATCAGATTG	GTGTGCACCT	TTCGCATCAT	AATCTTCCTA	TCTTGGGAGA	3300
CCCTCTCTAT	AATAGTAAAT	CAAAGACAAG	CGGGCTTATG	CTTCATGGCT	TCCGACTTTC	3360
CTTTACCCAC	CCACTTACTT	TAGAGAAGCT	AACCTTCAC	ACCTTTTCAA	ATACATTTGA	3420
AAAGAAATTA	AAAAAGAAATG	GATGATCGTG	TCATCCATTT	TTCCATATAA	AAAAGCAAGA	3480
CCACAAAGCC	TTGCTTTCTA	TCAACTCAAG	AATTATTTAG	CAATTTTTCG	GAAGTATTCA	3540
AGAGTACGAA	CAGGTGTGTC	ACTGTATGAC	ATTTCTGTTG	CGTACCATGA	TACAACCTTA	3600
ACCAATTGTT	TACCGTCAAC	GTCAAGAACT	TTAGTTTGAG	TTGCGTCAAA	CANTGAACCG	3660
TAAAGACATAC	CTACGATATC	TGAAGATAAG	ATTGGATCTT	CTGTGTAACC	GTATGATTCG	3720
TTTGAAGCTG	CTTTCTATAG	TGCGTTCACT	TCATCAACAG	TAACTGTTCT	TTCAAGAAGT	3780
GCTACCAATT	CAGTAACTGA	TCCAGTTGGA	GTTCGAACGC	GTTTTCGAGA	TCCGTCAAGT	3840
TTACCATTTA	ATTCTGGGAT	TACAAAGACG	ATAGCTTTTG	CAGCACCAGT	TGAGTTAGGA	3900
ACGATGTTTG	CAGCACCAGC	GCGAGCACGG	CGAAGGTCAC	CACCACGGTG	TGGTCCCTCA	3960
AGGATCATTTT	GGTCACCAGT	GTAAGCGTGG	ATAGTAGTCA	TCAATCCTTC	AACAACACCA	4020
AAGTTGTCTT	GAAGAGCTTT	AGCCATTGGA	GCCAGCAGT	TTGTAGTACA	TGAAGCACCT	4080
GAGATAACTG	TTTCAGTACC	GTCAAGAACG	TGCTGGTTAG	TGTTGAATAC	AACCTGTTTA	4140
ACCTGCTTTT	CACCAGGAGC	AGTGATAACA	ACTTTTTTAG	CTCCACCTTT	AAGGTGTTTT	4200
TCAGCTGCTT	CTTCTTAGC	AAAGAAACCA	GTAGCTTCAA	GAACGATTTT	TACACCTGCA	4260

		156	
GTAGCCAGT	CGATTGTTT	TGGATCACGT	TCAGCAGAAA
CTTTGATGAA	TTTACGGTTA	4320	
ACTTCAAAATC	CACCTCTTT	AACCTTCAACA	GTACCGTCGA
AACGACCTTG	AGTGTGTGCG	4380	
TATTTCGAACA	AGTGTGCAAG	CATAACTGGA	TCTGTAAAGT
CGTTGATGCG	TGTAACCTCA	4440	
ACACCTTTCTA	CGTTTGGAT	ACGACGGAAA	GCAAGACGAC
CGATACGTCC	GAACCGTTA	4500	
ATACCAACTT	TAATACCAT	TAGTGATTTC	CTCCTTATGA
AAATCATGAA	ATTTTATTG	4560	
TGAAAAGAGT	AACCTGAAATC	ACTACAAATC	ACCTTTCAC
AAACCTATT	TACAACTATT	4620	
TGAGTTGAAT	TGCAAGTATG	GCCATTGTTT	TTCTATGTTA
GTTTCTTTT	AAGACTGTAA	4680	
ACCAAGGAAT	CCCTTACTAT	TCATAGCATA	ACGATTCTAT
AGGATTCATT	TTACTAATCT	4740	
TACGCGCCG	GAAGTAGGCT	GAGACATAAC	CAGTAATAG
AGCGAAAAT	AGAGTTCCTA	4800	
AAACAGATAA	AAGATTTAAT	TTAAAAACCT	TAGTGATGGA
TGGGTAAAG	TGACTTACAA	4860	
TCGCACTTCG	CAAACTTCC	ACCCCTTGTC	CAACAAAAA
TGCCAGCAGC	AAGCGATGC	4920	
CTACAATCCA	GATAGCCTCG	TAAATAAAAA	TTCCCTTGAC
ATCACGATTC	TGATAACCA	4980	
CTGCTTTTAT	GACACCTATT	TCCCTGGAAC	GTTGCATGAT
ATTGAITGAA	ATAATGATAC	5040	
CAATCATAAC	CGCTGCTACC	ACAATAGCTT	GTGATGAAG
CACAATCAAT	AATCCCTGAA	5100	
TAAACGAAT	AAAGTAAATC	ACAAATACAA	GAATCTCTG
TTGAGAAAGC	ACAGTATACT	5160	
TCTTATTTTT	CTGTAATCT	TCTGTACTA	CTTTTGCTCG
TGATGGATCT	TTGAGTTCCA	5220	
AGATAAATA	AGATACAGCT	TTCTGTAATC	CAGCTCTTTC
CAAAATCGTT	TCCATTGAT	5280	
GAGACAGCAT	GAAGCTGTG	CTGTCTTCCA	TGTATCTTC
ATCATTCATT	ACACGTACAA	5340	
TCTTCGTTTG	AAATTGAGCA	ATCTTACTAG	TTTCGGCAGC
ACTTCTTACA	ATGCTGGCTG	5400	
AGACTGATTT	GCCAATAAGA	TCATTAGCTG	TCAAAATTTT
TCTGTCTGT	TCATTCCAA	5460	
TTTTTAGTAA	ACTGCTTGA	ATCGTTAATC	CCTGTTTATT
TGTATCAGTA	TAGAGGGATC	5520	
CAGCCAACAC	TTTGTCCGTC	TCATTATTAC	TAAACAGAGAT
ACTTGTATCA	TCATAAAGAC	5580	
TCACTACTTG	AGCATAAGAA	GGCATCGTTT	GACTCAGATC
CATTCTTTC	CCATCTATAG	5640	
TAATATTGTA	CATGTTTCATC	CCAAAAGGAC	TCTCCAATA
TTTAATAGCT	TCTTTCCAA	5700	
CTGTATCCGT	GATATATAGT	CAATTGAAAC	AAGACAGGA
TAAAAAGCC	TCGTAAAGG	5760	
TATTGCAACT	TGTAATACC	TTTTTGAGGT	GCTTTTGTAT
ATGAGCCCAT	GTTTCTCAA	5820	
TAGGATTGTA	CTCAGCGAG	TAGGAGGAA	GAGGTAAAG
TTTATGCCCA	AACCTTTCG	5880	
ATAAAAGTTC	TAGCTTCCCC	ATCTATGGA	ATCTTACATT
ATCCATAATA	ATAACOGATG	5940	
GTGTGTTTAA	TGTTGTGAAG	AGAAAAATCT	GAACCAAGC
TTCAAAAAAG	TCGCTGTCA	6000	
TCTCTCTTTC	GTAAGTCATT	GGAGCGATTA	ATTCAACATT
TGTTAGACCT	GCAACCAAG	6060	

157

AAATCCTCTG ATATCTTCTT CCAGATACCT TGCCCTCTAT TAATTGACCT TTTAATGAGC	6120
GACCATATTC TCGATAAAAA TAAGTATCGA ATCCGTGTTT GTCAATCTAA ACAGGTGCTA	6180
GGTGTCTTAA ACTATTAAAA TTCTTAAGAA ATAAGGCTAC TTTTCTGGG TCTTGTTCAT	6240
AGTAGGTGTG GTTCTTTTTT CGAGGTAGC CCATAGCTTT GAGCGTATAG TGGATGGTAG	6300
TTGCATGACA GCCAAATICA GAAGCTATTT CAGTCAAATA AGGCTCTGGA TTGTCAGTAA	6360
GATAGTTTTT AAGTCTATCT CTATCAACCT TTCTTGGTTT TATTCCTTT ACTTGGTGT	6420
TTAGCTCTCC TGTCTTCTCT TTAGCTTTA ACCAGCATA AATGGTATTA CGTGAGATTT	6480
GGAAACGTG TGTGCTTCT GTTATACCT CTGTTGCTC ACAATAAGAG AGAACTTTTT	6540
TACGAAATC TATTGAATAT GCATATAAAA GATTATACCA CATTTGTGTAC TATTTTGTGT	6600
TCATTTTACT ATATTTGAAG AGGCGTTTAA ACTATCTGAC ATAAAACCTG TTCTAGAGGA	6660
AAGACATCCT TTAATAAGTT AGTTTATTTT ACAACTTACA CATCAAGGTA GGTAAACCCC	6720
TTCAATGAAA AATCAAGACT CTTAGACTA TGGGTTAAAC TACCACTGGA GACGTAATCA	6780
ATCGCTAAAC CACGAAACG GCTAATAGTG GTCATATCAA TATTTCCAGA ACATTCATTC	6840
CGAGAACGTC CTGCAATTAG GGTAAATGGC TGTTCATCT GTTCCAATGA CMTATTATCC	6900
AACATGATAA TATCAGCACC CGCCGCCGCA GCTTCTTGG CAGCAGCAAG GCTTTCCACT	6960
TCCACCTCGA CCATTTTAC AAAAGGGGCA TAGGCACGG CTTGAGCAAT TGCCTTTTGA	7020
ACACTACCTA CTGCCGCAAT GTGATTGTCT TTTAGCAGGA TAGCATCTGA TAAATTAAAG	7080
CGATGATAT AGCCACGCC AACTCTCAG GCATATTTCT CAAAAGACG TAAATTAGGA	7140
GTAGTTTTTC GATATCAAA TACCTTAATG CAATCATGCG CTAAGGCTTC TACATAAGCA	7200
GCTGTCAATG AAGCAATCCC TGATAAATGT TGTAAAAAT TCAAGGCAAC GCGTCAACAT	7260
GTTAAGAGAC TTCTCACCGA GCCTATGATT TCTAAAAACA AATCGCACT AGTCAAACTA	7320
TCCCCATCCT TAAATTGATG AGGATCTGAG AAGGTCACTT CGGCATCAAA TAGGGTAAAA	7380
ACCCTTTGAA AAACGGTTAG CCCCCTAAA ACACCGCTT CCTTGGCAA AAGGACACCC	7440
TTGGCTTTGC CATGATGATC AAAAATGGCA TTGGTACTGT AATCTTCGGA ATGAACATCT	7500
TCTGCAAGC CTGCTTCAA TGTATCATCT ATTTGAAAAG GGGTTAAATC AGTTGAAATG	7560
ATTGACATCA C	7571

(2) INFORMATION FOR SEQ ID NO: 3:

- (i) SEQUENCE CHARACTERISTICS:  
 (A) LENGTH: 25385 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double

(D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 3:

TTTGCTAGTG GCTTAAATTC TTCAAGAAAA TCAGGCGTAT CTAAAGTCG TGTCGTTTTT	60
GTTCATCTA TATAAGACT TCCTGCTCCC CTACAACTA GAAAACGTGT CTGTGTPCCA	120
GCAAGAAGCT GATTAAATAG TTGAGATTGAT TTGCTGTGGA GCGGTAGCGT ATCTGGTGTGA	180
TAAGCACCAA ACGCTGAAAT AACAGCATCA AATCCAGTAA GATCATCTTT TGTCAACTCA	240
AATTAATCTT TTTTAATAAT AGACTCAGCT TGACTTTTGT TTTCAAGAAG AACAAATAGCC	300
GTTACTTCAT GTCCCTGGTT GACTGCTTCT TCAACAATTG CTTTCCCGCG TTGTCCATT	360
GCTGCAATAA CTGCTAGTTT CATTTTTTAT ACCTCTCTTG TTGTAATTAT TTTAGTTACA	420
GAAATTGTGA CACTCTTAAT AATCAATGTC AATAGTCTTG CTTAAATTAT ATCAAAATAT	480
TTCTACCAAG AALACTAAC ATGATTCTAG TGAATAAATA TCTTCTTTGT CAACAAATTT	540
ACTTTCTTGT TTTAAACATG CTATAATAAT CATAGCAAGA GATCTAAGTT GTCTGTTTTT	600
TTAAACGAG GTGATTATCA TGCCTAGATT CTATCCCAT CTCCCTACT ATCTGGTCAT	660
ATTATCTCTT TATTGGCCAC TTTATGAGTT GTTCTTACTA GTTCTTTCTG ACCCCCTTAC	720
ACTCAAGGGA CTCTATATAA ACAATCTTCT CTTCTTACCA CTTCTGGTAA TCTTGATTGT	780
ATGTTTACTC TATAGTACC GTTTCGGTTT CTCACTTTGA TGGTTAGTTG GTAACGGACT	840
GCTCTTTTAC TTTACTATCA TAACCTTTGG TGAGTTTATA CTAATTTACT TGCTAATCTA	900
TGAACACGTT GCTCTGGTCG GCATGGATTG TGGTATTAGC ATCAAGCATA TTCTACAAAA	960
AATGAAAAAC AAAAAACTTT CACAAAAATC TTGAAAAATC TCACAATCAT GCTATAATAA	1020
TCCATAGAGA CAAGTCACTT AGTCCCTTTC TACTAGAGAG TGCGTGGTTG CTGGAACGCG	1080
ATAGGAAGTC TAAACTGATA CTACTCTTGA GTTTTTFATG AAAACATAAA ACGGTGGCCA	1140
CGTTAGAGCC GATCAGAGGT GTCCCTCTCT TTTGAGGTAC ATAAATGAAG GTGGAACCA	1200
GTGTCAGCT CTTTTCGAGG ATGTGCAATT TTTTATTTAG GATACTAATT ATGGAGTTGC	1260
AAGAATTAGT GGAGCGCAGT TGGGCAATCC GACAAGCTTA TCACGAAGCTG GAAGTTAAGC	1320
ATCATGATTC CAAGTGGACG GTAGAAGAAG ACCTCTTGGC TTTATCTAAT GATATTGGAA	1380
ATTTCCACAG ACTGGTGATG ACAAAGCAAG GACGCTACTA TGATGAACAA CCCTACACAC	1440
TGGAACAAAA ACTTTCAGAA AATATCTGGT GGCTATTAGA ACTTTCTCAA CGTTTGATA	1500
TAGACATTCT GACGGAATG GAAAACTTCC TCTCTGATAA AGAAAAAGCA TTGAACGTTA	1560
GGACTTGGAA GTAGTCTGCT GATAAAAAAT CAATGCTTAG AACTATGAA ATAAATAAAA	1620

AGGAGAACAT	CATGATTAAC	ATTACTTTCC	CAGATGGCGC	TGTTCTGTGAA	TTGGAATCTG	1680
GCGTAACAAC	TTTGTGAATT	GCCCAATCTA	TCAGCAATTG	CCTPAGCTAAA	AAAGCCTTTGG	1740
CTGGTAAATT	CAACGGCAAA	CTCATCGACA	CTACTCGGCG	TATCACTGAA	GATGGAAGCA	1800
TCGAATTTGT	GACACCTGAT	CACGAAGATG	CCCTTCCAAT	CTTCCGTGAC	TCAGCAGCTC	1860
ACTTTGTTCC	CCAAGCAGCT	CGTCTGCTTT	TCCCAAGCAT	TCACTTGGGA	GTTCGTCCAG	1920
CCATCGAAGA	TGGTTTCTAC	TACGATACTG	ACAACACAGC	TGCTCAAATC	TCTAACGAAG	1980
ACCTTCTCTG	TATCGAAGAA	GAAATGCAAA	AAATCCTCAA	AGAAAACCTC	CCATCTATTC	2040
GTGAAGAAGT	GACTAAGAC	GAGGCACGTG	AAATCTTCAA	AAATGACCCT	TACAAGTTGG	2100
AATTGATTGA	AGAACTACTCA	GAAGACGAAG	GCGTTTGGAC	TATCTATCGT	CAGGCTGAAT	2160
ATGTAGACCT	CTGCCGTGGA	CCTCACGCTC	CATCAACAGG	TCGTATCCAA	ATCTTCCACC	2220
TTCTTCCATGT	AGCTGGTGGG	TACTGGCGTG	GAAACAGCGA	CAACGCTATG	ATGCAACGTA	2280
TCTACGGTAC	AGCTTGGTTT	GACAAGAAAG	ACTTGAAAAA	CTACCTTCAA	ATGCGTGAAG	2340
AAGCTAAGGA	ACGTGACCAC	CGTAAACTTG	GTAAGAGCT	TGACCTCTTT	ATGATTTTAC	2400
AAGAAGTGGG	ACAAGTTTTG	CCATCTCTGT	TGCCAAATGG	TGCGACTATC	CGTCGTGAAT	2460
TGGAACGCTA	CATCGTAAAC	AAAGAGTTGG	TTTCTGGCTA	CCAACAGTTC	TACACTCCAC	2520
CACTTGCTTC	TGTTGAGCTT	TACAAGACTT	CTGGTCACTG	GGATCATTAC	CAAGAAGACA	2580
TGTTCCCAAC	CATGGACATG	GGTGACGGGG	AAGAATTTGT	CCTTCGTCCA	ATGAACGTGC	2640
CGCACCATAT	CCAAGTTTTT	AAACACCATG	TTCACTCTTA	CCGTGAATTTG	CCAATCCGTA	2700
TGCGTGAAT	CGGTATGATG	CACCGTTACG	AAAAATCTGG	TGCCCTCACT	GGCCTTCAAC	2760
GTGTACGTGA	AATGTCACTC	AACGACGGTC	ACCTATTCTGT	TACTCCAGAA	CAAATCCAAG	2820
AAGAATTTCA	ACGTGCCCTT	CAGTTGATTA	TCGATGTTTA	TGAAGACTTC	AACTTGACTG	2880
ACTACCGCTT	CCGCTCTCTT	CTTCGTGACC	CTCAAGATAC	TCATAAGTAC	TTTGATTAAG	2940
ATGAGATGTG	GGAAAAATGCC	CAAACTATGC	TTCTGTGACG	TCCTGATGAA	ATGGGCGTGG	3000
ACTACTTTGA	AGCCGAAGGT	GAAGCAGCCT	TCTACGGACC	AAANTTGGAT	ATCCAGATTA	3060
AAACTGCCCC	TGGAAAAGAA	GAAACCTTTT	CTACTATCCA	ACTTGATTTT	TTTGTGCCAG	3120
AACGCTTCTGA	CCTCAAAATC	ATCGGAGCTG	ATGGCGAAGA	TCACCGTCCA	GTCAATGATCC	3180
ACCGTGGGTT	TATCTCAACT	ATGGAACGCT	TCACAGCTAT	CTTGATTGAG	AACTACAAGG	3240
GGGCTTCCCT	AACATGGCTG	GCACCACACC	AAGTAACCTT	CATCCCAAGTA	TCTAACGAAG	3300
AACACGTGGA	CTACGCTTGG	GAAATGGCCA	AGAAATCCCG	TGACCGCGGT	GTCCCTGCGAG	3360

160		
ACGTFAGATGA GCGCAATGAA AAAATGCACT TCAGATTCCT TGGTTCACAA ACCAGCAAGA	3420	
TTCTCTTACCA ATTAATTGTT GGAGACAAAG AAATGGAAGA CGAAACAGTC AACGTTGGTC	3480	
GCTACGGCCA AAAAGAAACA CAACTGTCT CAGTTGATAA TTTTGTCAA GCTATCCTAG	3540	
CTGATATCGC CAACAATCA CGCGTTGAGA AATAAGATC TAGCATAAAA GCCTCCAATC	3600	
TGGAGCGCTT TTCTCATCTA TTTTACTCA AGGACTAAGT TCACCTGACC AAATGAATC	3660	
CGCACTGTGC TTCTTTTCC GACCTCAGAC TCGATACGAA TCTGCTGCC CAGTTCTTCA	3720	
GAAATTTTCT TAGATAGATA AAGGCCAAGT CCAGAGGACT GCTGGGTCAA ACGGCCATTG	3780	
TATCTTGAAA AGCCACGCTC AAATACTCG AGGACATCAC TGTTTTTAT CCCGATTCCC	3840	
GTATCTTTGA TACAAGCTC TTGCTCATCC ATATAAATCT CCAGACCACC TTCCTGGTG	3900	
TACTTGAGAC TGTTTGAGAT GATTTGCTCA ATAACCACTA GCAGCCACTT TTTATCCGTC	3960	
ACGATTTCTT TATCAAGGTC ATGAGATTG ACATTTAAGC CTTTTGAAT AAAGAAAAGA	4020	
CGATATTTAC GAATPATTC CTGACCAAG TCCTCAATT GAACTGCTT TAAGACCAA	4080	
TCATCATGGA AACTTTCTAA ACGCAGGTAC TGTAAACTA GGTGGTATA GGAGTCGATT	4140	
TTGAAAATTT CCTGTTCTAG CTGCTGCTTC AGTTGGCGGT CGACCACTTC TCCACTAAG	4200	
AGTTGACTGG CTGCAATGG GGTCTTTATC TGATGGACCC ACAAGGTATA GTAATCCAGC	4260	
AAATCCGTCA GTTTCTTTTC TGCTTTTGAC CTCTGCTGAT AGAGTCCAT CTCACGGGCT	4320	
TCTAATTTT CTGCTAAGC TATTTCCAAA CGAGACTTGG CTTCCCTCTC TCCATAGAGA	4380	
AGTTCTGGC GATAGACCTG CGTTCCACC AATATGTCCC AAGTGAAAA TAATATGGTT	4440	
ACAAAGCAAC ACAAGAAGAA AAGTAGAGG AAGTAAATC CTAGACTGCC AAATAAAAC	4500	
TGAAGAAGTA AGACAAGAAA TGCCAAAGAA AGCAGATAGA TAAAAAGAC ACTACGGGAG	4560	
CGCAGATAGG CTAGAAAAAA TTGTTTCCAA TCAAGCATGC TTCAATCGT ACCCTATTC	4620	
TTTCTTGGTC TCGATAAATC CTACCAATCC CTGCTCCTCC AACTTTTAC GCAAAACGAC	4680	
CACATTGACA GAGAGGGTAT TATCATCAAT GAAAAAGTCA CTGTTCCAAA GTTCCCGCAT	4740	
CAGGTCGTCA CGTGCTACGA TGTTCCTGCT ATGCTCAAAT AACACCGGTA AAATCTGGAA	4800	
TTGATTTCTG GTCAAAATCA AGACTTGCC TTGATAATGT AAATCCATGG ATTTGGTATT	4860	
GAGGATAACA CCAGCATATT CCAGCAAACT CTCATCAGC CCAAACTCAT AGGAACGACG	4920	
CAACAAGCCC TGAACCTTAG CTAAAAGAAC CTGCTGTGTA AAAGGCTTGG TCACAAGTCC	4980	
ATCCGCCCCC ATATTGATTG CCATGACAAT ATCCATAGCC TGGTCTCTCG AAGAAAGAAA	5040	
CATGATAGGT ACCTTGAAA TCTTGGGAT TTCTGACAC CAGTGATAAC CATTAACAA	5100	
GGGCAAACCA ATATCCATGA GGACCAGATG AGGTTCCGAC TGAACAAATA GACTCAAAAC	5160	

TTCCATAAAG	TTCTTACCA	GGACCACCTC	AAATCCCAT	TCAGAGAGCA	TTTCCCAAT	5220
CTGTGACGA	ATGACCTGAT	CATCTTCTAT	TAATAAAATC	TTGTGCAAGC	GCTTCTCCTT	5280
TTCCATTATT	ATAACAGATT	TTTCCATGCT	AGATGGTCIG	AAACGAAAT	TGAAATAGCC	5340
TGTFTTTAGC	CAGTACAAC	AGGCTATGCT	ACTAGCTAAT	TTGAGGAAA	TTTGCTAAGA	5400
TAAATAAATA	GAAAGAGAGT	CTTATGGCCA	ATATTTTGA	CTATCTGAAA	GATGTCCGAT	5460
ATGATTCTTA	TTACGACCTT	CCCTTGAATG	AGTTAGACAT	TCTAACCTTA	ATAGAAATCA	5520
CCTACCTCTC	CTTTGATAAT	CTGCTCTCCA	CACITCCTCA	ACGCTCTTTA	GATCTAGCAC	5580
CTCAGGTCC	AAGAGATCCC	ACCATGCTTA	CPAGCAAAAA	TGCGCTTCAA	TTATTAGATG	5640
AAATGGCTCA	ACACAAGGCG	TTCAAAATTT	GCAAACTCTC	CCATTTTATC	AACGACATCG	5700
ACCGTGAATC	GCAAAAGCAA	TTTGGGGCTA	TGACTTATCG	TGTCAGCCTC	GATACCTATC	5760
TGATTTGCTT	TCGTGGGACA	GATGACAGTA	TCATTTGGCTG	GAAGGAAGAT	TTCCACCTGA	5820
CCATATATGA	GGAAATTCCT	GCTCAAAAGC	ACGCCCTTCG	CTATTTAAAG	AACCTTTTTC	5880
CCCATCATCC	TAAGCAAAAG	GTATTTCTAG	CTGGGCATTC	CAAGGGAGGA	AATCTCGCTA	5940
TCTATGCTGC	TAGCCAAAT	GAGCAAGTT	TGCAAAATCA	GATCACAACA	GTTTATACAT	6000
TTGATGCACC	TGGTCTCCAT	CAAGAAATGA	CACAGACTGC	GGGTATCAA	AGGATTAATG	6060
ATAGAGACAA	GATATTCATT	CCACAAGGTT	CCATTATCGG	TATGATGCTG	GAATTTCTTG	6120
CTCACCAAA	CATCGTTTCA	AGTACTGCC	TGGGTGGCAT	CGCCAGCAC	GATACCTTTA	6180
GTTCGCAGAT	TCAGGACAA	CACCTCGTCC	AACTGGATTA	GACCAACAGT	GATAGCCAGC	6240
AACTAGACAC	AACCTTTAAA	GAATGGGTGG	CCACAGTCCC	TGACGAAGAA	CTTCAGCTCT	6300
ACTTCGACCT	CTCTTTGGC	ACTATTTCTG	ATGCTGGTAT	TAGCTCTATC	AATGACTTTG	6360
CTTCTCTAAA	GGCCTCTGAA	TACATTCATC	ATCTCTTTGT	CCAAAGCTAA	TCCCTCACTC	6420
CAGAAGAAAG	AGAAACCTTG	GCTCGCCCTA	CCAGTTATT	GATPGATACT	CGTTACCAGG	6480
CATGAAATAA	TAGATAATAC	CTTTGAAAT	TAAATGTATA	CAAAACAAAA	GACCTAGAAAT	6540
ACATATCTTC	ATGTGCATTC	TAACTCTTTT	TAAATAGAAAT	CTAATAGTCA	ATAAAAAATCA	6600
AAGACATATG	AGAGATTAATG	GGCCTTGGAA	CGTCCCTCTC	GCTTCAACAA	AATGACCCCA	6660
TTATAGATTA	AAAAGATGCC	ACTTAGAAAA	AOCAAAAAAG	GAACTAAGAC	AAAGGCCAAAT	6720
ATATAAAAAG	CTAACTGAAC	ATTCTCGTAT	CCATTTTTAT	AAAAAGGTA	GGATAGATTA	6780
AAATAACTTG	AAATGAGGGA	TAATAAAAAAT	AATACTGGAT	TCCAACAATC	TCTATTATCC	6840
TTCCAAAAATG	ACACTATAAA	GGCTAATACA	ATTCCTATAA	CGAGATACAT	TTCTTACTCC	6900

162						
TTTAATAGCT	ACATTTTATC	ATAATTTATC	AAAGAAAAAA	GAGGGCATTT	ATCCCTCTTA	6960
ATCCTTCACT	TGACTCTCTG	CATCGGCCAC	GACTTTTCTT	AGACTGGTCT	GACCAAGTTC	7020
TGCCCTCATA	GTCAACTGAA	TTCTCTCCAA	TTTTTGATCC	AAAACATCAT	GAATATGAGC	7080
TCTACAGGG	CAATTTGGAT	TCGGATTGTC	ATGGAAACTG	AAGAGTTGAC	CTGTCTTACC	7140
AAGACATTGC	ACCGCTGAT	AAACATCTAA	AGACTAATA	TCCTTAAGGT	CCTTGACAA	7200
CTCTGTTCCG	CCCGTTCCAC	GGCTACTGTA	AATCAGCTCT	GCCTTCTTCA	ACTGGGACAA	7260
GATCTTTCCT	ATAATGACAG	GATTGACCCC	GACACTAGCA	GCCAGAAAA	CACTGGTCAC	7320
CTTGCTTTCC	TTCCCTTCGA	GGGCAATGAT	TATCAGCATA	TGAGTCGCAA	TGGTAAATCT	7380
ACTTGGAATT	TGCATCTCTC	TCTCTTTTT	ACGAGGTACT	CCTGCCTCTA	CTCTCTTTTT	7440
TCTATTATTA	TACCTTTTTT	AGTTGTAATG	TCAATGGTTA	CCACTTTTCA	ACCACTCGTC	7500
TAACATCCGA	TGCGAGCCCT	CTTTCTGAGC	CAATTCTCTC	AAAAATTCCT	GATGATGAGT	7560
ATGGTGGATC	CCATTGACCA	GACTTTTATA	GTAAACCTCA	AAATAGGGAA	GTCTCAGTTC	7620
TTTAGCCAGC	TGCAATTTCAG	CTGCTACATC	GTAGTCTACC	CGTCGGAAAT	CCATATCTAC	7680
CAGGCTTTTG	TCATCAAATC	CCAAAATCAT	ATACTGGGCC	CGCAAGTCTT	TCCGTAGCTG	7740
AGCGTCCAAA	AAGAAGGTT	GGCCAATCGA	ACCCGGAATT	ACAAATCAAT	GCCCACCAGT	7800
CCCGTAACGA	AGCAACTGCT	GGTGAATATG	TCCATAAACA	GCAATATCAC	AGGGAGGATG	7860
AGTCACCAAG	CGGTCAAATC	CCTCTTGTTT	GCCAGTATGA	ATCAACTCTC	GCCCCAGTTT	7920
CTTATCAGGC	AGATGATGGC	TAATTCCTCC	CGTCAAATCC	CCAAACTGAC	GATGAATTTG	7980
AAGAGGTTGA	TTGTGAGACA	CTTCAATTTT	TTCTAGGGAA	ATTTCTCTTA	AAACATACTG	8040
GCACGTGGCG	AAGAGATAGC	GTGACTGGG	GCGAGTACTG	TCCAATTCCT	TACGGACACC	8100
ATGCCAAGAA	CTGTCTTCCC	AGTTTCCCAA	AACTCTAGCC	GTAAATCGTA	GTTGATCCAA	8160
CAAGTCCAAA	ATCTTTCTAC	GCCCTGTCCC	TGGCATGAGA	ATATCTCCCA	AAAGCCAGTA	8220
TTTATCCACT	CCTATCTGCC	GAGCATCTGC	CAAAACAGCC	TCCAAGCGGG	TGGTATTTCC	8280
ATGAATATCT	GAAGAAGAG	CTATTTTCGT	CATATCCATC	TCTCTGTTT	TTCTCTTGCA	8340
ATAAGTATAA	CATAAAAAAT	CACAGCTAGA	GAATCTAGC	TTTTTTTGAT	ATACTAGATA	8400
AAGATATTAG	ACAAAGAGAA	ACGAATGACC	CCAAACAAG	AAGACTATCT	AAAATGTAT	8460
TATGAJATTG	GCAATAGACCT	GCATAAGATT	ACCAACAAGG	AAATTCGGCG	TGCGATGCAA	8520
GTCTCTCCCC	CTGCCGTAAC	TGAAATGATC	AAAGCATGA	AAAGTGAAAA	TCTCATCCTA	8580
AAGGACAAAG	AATGTGGCTA	TCTACTGACT	GACCTCGGTC	TCAAACTGGT	CTCTGAGCTC	8640
TATCGTAAGC	ACCGCTTGAT	TGAAGTTTTT	CTAGTTCATC	ATTTAGACTA	TACAAGTGAC	8700



CAGATTCAAG	AGGAAGCTGA	GGTCTTCCAA	CACACTGTCT	CTGACCTGTT	CGTGGAAAGA	8760
CTAGATAAAC	TGCTAGGTTT	CCCTAAAACC	TGCCCCCAAG	GGGGAACAT	TCTCGCCAAG	8820
GGAGAAGTAC	TCGTTGAAAT	CAATAACCTC	CCACTAGCTG	ATATCAAGGA	AGCTGGCGCC	8880
TACCGCCTGA	CTCGGGTGCA	CGATAGTTTT	GACATTCTCC	ATTATCTGGA	CAAGCACTCA	8940
CTTCACATCG	GTGACCAGCT	CCAAGTCAAG	CAGTTTGATG	GCTTCAAGCA	TACCTTCACT	9000
ATCCTCAGTA	ACGACGAGGA	TTTACAAGTG	AATATGGACA	TTGCAAAACA	ACTCTATGTC	9060
GAGAAATCA	ACTAATTTCT	CAAGTCCOCT	ACCAACCCTG	AAAGTTTTAT	TTTGCTCTTT	9120
TGTCAACTGT	AGTGGGTTGA	AGTCAGCTAA	GCTCGAGAAA	GCACAAATTT	TGTCCTTTCT	9180
TTTTTGATAT	TCAGAGCGAT	AAAAATCCGT	TTTTTGAAGT	TTTCAAAAGT	CCGAAAACCA	9240
AAGGCATTGC	GCTTGATAG	TTTGATGAGA	TTATTGGTCG	CTTCCAGTTT	GGCATTAGAA	9300
TAGTGTAGTT	GAAGGCGCTT	GACAATCTTT	TCTTTATCTT	TGAGGAAGGT	TTTAAAGACA	9360
GTCAGAAAA	TAGGATGAAC	CTGCTTTAGA	TTGTCTTCAA	TGAGTCCGAA	AAATTTCTCC	9420
GGTTCTTTAT	TCTGAAAGTG	AAACAGCAAG	AGTTGATAGA	GCTGATAGTG	GTGTTTCAAG	9480
TCTTTGTAAT	AGCTCAAAAG	CTTGCTTAAA	ATCTCTTTAT	TGGTTAAGTG	CATACGAAAA	9540
GTAGGACGAT	AAAATCGCTT	ATCACTCAGT	TTACGGCTAT	CTGTGTTGAT	GAGCTTCCAG	9600
TAGCGCTTGA	TAGCCTTGTA	TTCTATGGAT	TTTCGATCCA	ATTGGTTTCA	AATTTGAACA	9660
CGCACACGAC	TCATAGCAGC	GCTAAGATGT	TGTACAATGT	GAAAGCGATC	CAACACGATT	9720
TTAGCATTCG	GGAGTGAAAC	AGTCTGGGAG	ACTGTTTCAG	CCTGAGCCTA	GAAATTTGAA	9780
AGCGAAGCTG	TTTAGCCAAAG	TCATAGTAAG	GACTAAACAT	ATCCATCGTA	ATGATTTTCA	9840
CTTGACAACG	AACGGCTCTA	TCGTAGCGAA	GAAAGTGATT	TCGGATGACA	GCTTGTGTTT	9900
TGCGTTTCAAG	AACAGTGATA	ATATTAAAGT	TATCAAAATC	TTGCGCAATG	AAACTCATCT	9960
TTCCCTTAGT	GAAGGCATAC	TCATCCCAGG	ACATAATCTT	TGGAGCCGGA	GAAAAATCAT	10020
GCTCAAAAGT	AAAGTCATTG	AGCTTGGCGA	TGACAGTTGA	AGTTGAAATG	GCCAGCTGAT	10080
GGGCAATATC	AGTCATAGAA	ATTTTTTCAA	TTAACTTTTG	AGCAATTTTT	TGCTGTATGA	10140
TACGAGGGAT	TTGGTGATTT	TTCTTTTACCA	GGGAGTCTTC	AGCAACCATC	ATTTTTGAAC	10200
AGTGATAGCA	CTTGAAACGA	CGCTTTCTAA	GGAGAATTCT	AGAAGGCATA	CCAGTCGTTT	10260
CAAGATAAGG	AATTTTAGAA	GGTTTTTGAA	AGTCATATTT	CTTCAATTGG	TTTCCGCACCT	10320
CAGGGCAAGA	TGGGGCGCTG	TAGTCCAGTT	TGGCGATGAT	TTCTCTTGTT	GTATCCTTTAT	10380
TGATGATGTC	TAAATCTGG	ATATTAGGGT	CTTTAATGTC	TAGTAAATTT	GTGATAAAAT	10440

164

GTAATTGTC CATATGATTC TTCTAATGA GTTGTTTGT GCTTTTCAT TATAGTCAT	10500
ATGGGACTTT TTTTCTACAA TAAATAGGC TCCATAATPAT CTATAGTGA TTTACCCACT	10560
ACAAATATTA TAGAACCGTA AAAATAGAGT GAGATAGCAG GTTTCTAAGC CTGCTATCTT	10620
TTTGTGATGA CATTCAGGCT GATACGAAT CATAAGAGGT CTGAAACTAC TTTGAGAGTA	10680
GCTGTTCTTA TAAATATAG TAGATTGAAA TAAGATGIGA ACAACTCTAT CAGGAAGTC	10740
AAATTAAATT ATAGATTAT TTTAGCAGTC AAGGTGTACT GTTATAGATT CAATATATTA	10800
TATGACTATT AACCTTGCT TCTCTAAAA TTGACTTTCT TGTTTCTTA TCTGTCCAC	10860
TCGAAACAAG TACTGTAAAG ATTTGATTAT TTTTGAAAGT ACTTTTAATA TACTGTATAT	10920
AGTTAAAAAA GATTGAAAC TAAATTCCAA ATTGAAAAA GACTTGAAT ACTAAAAA	10980
AAAAAGTATA CTCTAATTGA AAACGTAAC AAAACTAATT TAGAAGATGA AATATAGAT	11040
ATTTCTCTCT TAAAGTTTT TGGTGAAAG AGATGTAGAA AGGAGATTTA GCCAAGAGT	11100
CTATTAGTGC TAGAATAATA GATTAGAAAT ATTTAGAAA AACGAAGTA GCAGCTTATA	11160
AAATCAAGTC CCCAANTAGA TTTACTACTAG TATCTTTTC AAAAAATAA GGGCGACTTC	11220
CTTCATGAAT ATCAATTCA TCTATAAGGA AGGTAGCTAA TTGAACTAAC TTATTTATTC	11280
TGTTGTGCGC TAGAAAAATC AGACCTCCTT GTGAAGATTG AGGAGATACT TAATGAAAA	11340
CAAGAAGAAA ACTAGCAAGC TAGTAGCAGA TTGCCAAAA CACCGCTTTG AGGTGTAGA	11400
TAAGACTGAC CTATATAATC CAAGGTGAG CGACTGTGT TTGAAGAGAT TTTCAAGAG	11460
TATAGGCTAG AGAGTAGTGT TTTTATGCTC TTCTAGTAGA AAATGCTAGA CAGAAGAAAG	11520
GGGAACCTGG ATAGGAAAA TAGATTGAGA AAGGAGGTTA GAAGAGATGA TTATTACAAA	11580
AAATTAGCCG TTAGGAACCT ATGTGGAGT AAATCCACAT TTTGCAACAT TAATAGATTT	11640
TCTAGAAAA ACAGGACTAG AAAATTAAAC AGAAGGTTG ATTGCTATCG ATGGTAATCG	11700
ATTGTTTGGG AATTGCTTTA CTATCTAGC AGATGGTCAA GCAGGGGCTT TCTTTGAAAC	11760
CCACCAAAAA TATTTGATA TTCAATTAGT TTTGGAAGC GAAGAAGCCA TGGCTGTAC	11820
ATCGCCGGAA AATGTAAGCG TTACCCAAGA ATATGATGAA GAGAAAGATA TTGAATTATA	11880
CACAGGAAJA GTGGAACAGT TGTTTATTT GAGAGCTGGC GAATGCCCTA TCACTTTTC	11940
AGAAGATTGA CATCAACCCA AGGTTCGTAT AAATGATGAA CCTGTGAAAA AAGTTGTCAT	12000
TAAAGTTGCG ATTTCTTAAT GTAGAAAGAG AAGAACGATG AAAAAAATGA GAAAGTTTTT	12060
ATGTCTAGCT GGAATTGCGC TAGCGGCTGT TGCCTTGATA GCTTGTTCAG GAAAAAAGA	12120
AGCTACAAC AGTACTGAAC CACCAACAGA ATTATCTGGT GAGATTACAA TGTGGCACTC	12180
CTTACTCAA GGACCCGTT TAGAAAGTAT TCAAAAATCA GCAGATGCTT TCATGCAAAA	12240

GCATCCAAAA	ACGAAAAATCA	AGATGAAAAAC	ATTTCTTGG	AATGACTTCT	ATACTAAATG	12300
GACTACAGGT	TTAGCAAAATG	GAATGTGCC	AGATATCAGT	ACAGCTCTTC	CTAACCAAGT	12360
AATGGAAATG	GTCAACTCAG	ATGCTTTGGT	TCGGCTAAAT	GATTCTATCA	AGCGTATTGG	12420
ACAAAGATAA	TTTAACGAAA	CTGCCTTAAA	TGAAGCAAAA	ATCOGAGATG	ATTACTACTC	12480
TGTTCTCTCT	TATTCACATG	CACAAGTCAT	GTGGGTTAGA	ACAGATTTGT	TAAAGAAGCA	12540
TAATATTGAG	GTTCCTAAAA	CTTGGGATCA	ACTCTATGAA	GCTTCTAAAA	AATTGAAAGA	12600
AGCTGGAGTT	TATGGCTTGT	CTGTTCCGTT	TGGAACAAAT	GACTTAATGG	CAACACGTTT	12660
CTTGAACCTC	TACGTACGTA	TTGGTGGAGG	AAGCTCTCTA	ACAAAGATC	TTAAAGCAGA	12720
CTTGACAGGC	CAACTGTCTC	AAGATGGTAT	TAAATACTGG	GTTAAATTGT	ATAAGAAAT	12780
CTCACCTCAA	GATTCTTTGA	ACTTTAATGT	CCTTCAACAA	GCTACCTTGT	TCTATCAAGG	12840
AAAAACAGCA	TTTGACTTTA	ACTCTGGCTT	CCATATCOGA	GGAATTAAAT	CCAACAGTCC	12900
TCAATTGATT	GATTGAGTTG	ATGCTTATCC	TATTCCAAAA	ATCAAAGAGT	CTGATAAAGA	12960
CCAAGGAATT	GAAACCTCAA	ACATTTCCAAT	GTTTGTTTGG	AAAAATTCAA	AACATCCAGA	13020
AGTTGCTAAA	GCATTCCTTAG	AAGCACTTTA	TAATGAAGAA	GACTACGTTA	AATTCCTTGA	13080
TTCAACTCCA	GTAGGTATGT	TGCCAATCAT	TAAGGGGATT	AGCGATTCTG	CAGCCTATFAA	13140
AGAAAAATGA	ACTCGTAAGA	AATTTAAACA	TGCTGAAGAA	GTAAATTAAGT	AAGCTGTATA	13200
AAAAGGTACT	GCTATTGGTT	ATGAAAATGG	GCCAAGTGTG	CAAGCTGGTA	TGTTGACTAA	13260
CCAACACATT	ATTGAACAAA	TGTTCCAAGA	TATCATTTACA	AATGGAACAG	ATCCTATGAA	13320
AGCAGCAAAA	GAAGCAGAAA	AACAAATAAA	TGATTTATTT	GAGGCTGTTT	AGTAGATGTA	13380
AAAGACTAGA	AAATAGGTGG	GATAGTGAGC	TGAAAAGCTC	TAGCCCAATC	TTGTAAAAGA	13440
AGGGAGAAGG	AGAAATGTTA	AAGAACGTAA	TTTAACCTGC	TGGATATTGG	TTTTGCCAGC	13500
TATGATTATC	GTAGGATTAC	TCTTTGTTTA	TCCGTTTCTC	TCGAGTATTT	TTTATAGCTT	13560
TACCAATAAG	CATTGATTAT	TGCCATAATTA	TAAATTTGTT	GGTTTGCTCA	ACTATAAAGC	13620
TGTCTATACA	GATCCCAACT	TCTTTAATGC	GTTCTTTAAT	TCAANTAAAT	GGACCGTTTT	13680
CTCATTAGTT	GGTCAAGTTT	TAGTAGGGTT	TGTATTGGCT	TTAGCTCTTC	ACAGAGTAGG	13740
CCACTTCAAG	AAATTATATA	GGACATTATT	GATTGTTTCT	TGGGCATTTT	CTACCATGTT	13800
TATTGCCCTC	TCTTGGCAGT	GGATTCTAAA	CGGGGTTTAT	GGCTACTTAC	CTAATCTAAT	13860
CGTAAATTA	GGTTAATGG	AACAATACACC	TGCATTTTTC	ACAGATAGTA	CATGGGCATT	13920
CCTATGTTTG	GTGTTTATCA	ACATTTGGTT	TGGAGCACCA	ATGATTATGG	TTAATGTGCT	13980

166						
TTCCAGCTTTTG	CAACACGTAC	CAGAAGAACA	ATTTCAGGCT	GCTAAGATAG	ATGGTGCTTC	14040
AAGTTGGCAG	GTGTTCAAGT	TTATCCTCTT	TCCACATATT	AAAGTGGTTG	TAGGACTTCT	14100
AGTTGTTTGG	AGAACTGTAT	GGATCTTTAA	TAACTTTGAC	ATTATCTACC	TCATTACTGG	14160
TGGTGGACCA	GCCAACTGTA	CAACGACGCT	TCCAAATTTT	GCTTACAACC	TGGGCTGGGG	14220
AACATAAATG	TTGGGTGCTG	CTTCAGCAGT	TACAGTACTG	CTCTTTATCT	TCTTGGTGGC	14280
GATTTGCTTT	ATCTACTTTG	CTATCATCAG	TAAGTGGGAA	AAGGAGGGTA	GAAAATAATG	14340
AAGAAGAAAT	CCAGTATTTA	TTTAGATATT	CTCTCACATG	TACTTTTAGT	TGGTCCGACC	14400
ATCGTTGCGAG	TTTTCCCATT	GGTATGGATT	ATCATATCTT	CTGTCAAAGG	GAAAGGGGAA	14460
TTAACTCAGT	ATCCAACACG	ATTTTGGCCT	GAACAGTTTA	CATTAGATTA	TTTCACTCAT	14520
GTTATCAACG	ATTTGCACCT	CATTGATAAC	ATTGGAALAA	GTTTAATCAT	TGCCTTGGCT	14580
ACAACCCCTTA	TTGCGATTAT	TATTTCTGCT	ATGGCAGCCT	ATGGTATTGT	TCGATTCTTT	14640
CCTAAATATG	GACCAATCAT	GTGAGACTTA	CTCGTCATTA	CCTACATTTT	CCCACCAATT	14700
TTGTTAGCAA	TTCCCTATTC	AATTGCCATT	GCTAAAGTTG	GGTTAACAAA	TAGTTTATTT	14760
GGCTTGATGA	TGGTTTATCT	ATCTTTTAGT	GTTCCATATG	CAGTTTGGCT	CTTAGTTGGA	14820
TTTTTTCCAAA	CAGTTCCAAT	TGGAATTGAA	GAAGCGGCTA	GAATTGATGG	TGCAAAATAA	14880
TTTGTACGTT	TTTATAAAGT	TGTGCTACCG	ATTGTAGCAC	CAGGTATTGT	AGCAACAGCT	14940
ATTATACAT	TTATCAATGC	TTGGAATGAA	TTCTGTATAT	CCTTGATTTT	GATTAACAAT	15000
ACAGGAAAGA	TGACAGTAGC	AGTAGCCCTT	CGTTCACCTA	ATGGTTCAGA	AATACTAGAC	15060
TGGGGAGATA	TGATGGCAGC	GTCTGTTATT	GTAGTCTCTC	CATCAATTAT	TTTCTCTCTC	15120
ATCATCCAAA	ATAAGATTGC	AAGTGGATTA	TCAGAAGGAT	CTGTGAAGTA	GACGAAAGAA	15180
GGAAAAAAAT	GAATAAAAGA	GGTCTTTTAT	CAAACTAGAG	AATTTCCGTT	GTAGGCATTA	15240
GTCTTTTAAAT	GGGAGTCCCC	ACTTTGATTTC	ATGCGAATGA	ATTAACTAT	GOTCAACTGT	15300
CCATATCTCC	TATTTTTCAA	GGAGGTTCAT	ATCAACTGAA	CANTAAGAGT	ATAGATATCA	15360
GCTCTTTGTT	ATTAGATAAA	TTGTCTGGAG	AGAGTCAGAC	AGTAGTAANT	AAATTAAGAG	15420
CAGATAAACC	AAACTCTCTT	CAAGCTTTGT	TTGGCTTATC	TAATAGTAAA	GCAGGCTTTA	15480
AAATAATTA	CTTTTCAATT	TTCAATGAGAG	ATTCTGGTGA	GATAGGTGTA	GAATAAGAG	15540
ACGCCCAAAA	GGGAATAAAT	TATTTATTTT	CCAGACCAGC	TTCATTATGG	GGAANAACATA	15600
AAGGACAGGC	AGTTGAAAAAT	ACACTAGTAT	TTGTATCTGA	TTCTAAAGAT	AAAACATACA	15660
CAATGTATGT	TAATGGAATA	GAAGTGTCTC	CTGAAACAGT	TCGATACATTT	TTGCCAATTT	15720
CAAAATATAA	TGGTATAGAT	AAGGCAACAC	TAGGAGCTGT	TAATCGTGAA	GTAAGGAAC	15780

ATTACCTCGC AAAAGCAAGT ATTGATGAAA TCAGTCTATT TAACAAAGCA ATTAGTGATC 15840  
 AGGAAGATTTC AACTTATCCC TTGTCAAAATC CATTTTCAGTT AATTTTCCAA TCAGGAGATT 15900  
 CTACTCAAGC TAACATTTTT AGAATACCGA CACTATATAC ATTAACTAGT GGAAGACTTC 15960  
 TATCAAGTAT TGATGCACCT TATGGTGGGA CTCATGATTC TAAAGTAAG ATTAATATTG 16020  
 CCACCTCTTA TAGTGATGAT AATGGGAAAA CGTGGAGTGA GCCAATTTTT GCTATGAAGT 16080  
 TTAATGACTA TGAGGAGCAG TTAGTTTAACT GGCCACGAGA TAATAAATTA AAGAATAGTC 16140  
 AAATTAGTGG AAGTGCTTCA TTCATAGATT CATCCATTGT TGAAGATAAA AAATCTGGGA 16200  
 AAACGATATT ACTAGCTGAT GTTATGCCTG CGGGTATTGG AATAATTAAT GCMAATAAAG 16260  
 CCGACTCAGG TTTTAAAGAA ATAAATGCTC ATTATTATT AAAACTAAAG AAGAATGGAG 16320  
 ATAACGATTT CCGTTATACA GTTAGAGAAA ATGCTCTCTT TTATAATGAA ACAACTAATA 16380  
 AACTTCAJAA TTATACTATA AATGATAAGT ATGAAGTTTT GAGGGAGGA AAGTCTTTAA 16440  
 CAGTCGAACA ATATTCCGTT GATTTTGATA GTGGCTCTTT AAGCAAAAG CATTAATGAA 16500  
 AACAGGTCC TATGAATGTT TTCTACAAAG ATTCGTTATT TAAAGTGACT CCTACTAATT 16560  
 ATATAGCAAT GACAACTAGT CAGAATAGAG GAGAGAGTTC GGAACAAATT AACTTGTTGC 16620  
 CTCCTGTTCTT AGGAGAAAA CATTAATGAA CTACTTATG TCCCGGACAA GGTTTAGCAT 16680  
 TAAAAATCAAG TAACAGATTG ATTTTTCGCA CATATACTAG TGGAGAACTA ACCTATCTCA 16740  
 TTTCTGATGA TAGTGGTCAA ACATGGAAGA AATCCTCAGC TTCAATTCCG TTTAAAAATG 16800  
 CAACAGCAGA AGCACAATG GTTGAAGTGA GAGATGGTGT GATTAGAACA TTCTTTAGAA 16860  
 CCACTACAGG TAAGATAGCT TATATGACTA GTAGAGATTG TGGAGAAACA TGGTGAAG 16920  
 TTTGTTATAT TGATGGAATC CAACAACTT CATATGCCAC ACAAGTATCT GCAMTTAAAT 16980  
 ACTCTCAATT AATTGATGGA AAGAAGCAG TCATTTTGAG TACACCAAAAT TCTAGAAGTG 17040  
 GCCGCAAGGG AGGCCAATTA GTTGTCCGTT TAGTCAATTA AGAAGATGAT AGTATTGATT 17100  
 GGAATATACA CTATGATATT GATTTGCCCT CGTATGGTA TGCCATATCT CGGATTAACAG 17160  
 AATTGCCAAA TCATCACATA GGTGTACTGT TTGAAAAATA TGATTCTGTG TCGAGAAATG 17220  
 AATTGCATTT AAGCAATGTA GTTCAGTATA TAGATTGGA AATTAATGAT TTAAACAAAT 17280  
 AAAGGAGAAA AACATGGTTA AATACGGTGT TGTGGGAACA GGGTATTTTG GAGCTGAATT 17340  
 GGCTCGCTAC ATGCAJAJGA ATGATGGAGC AGAGATTACT CTTCCTATAG ATCCAGATAA 17400  
 TGCAGAGCGC AATGCAAGAG AATGGGAGC AAAAGTAGCA AGTTCCTTAG ATGAGTTGCT 17460  
 TTCTACGCAT GAACTAGATT GTGTTATCGT CGCAACTCCA AATAATCTTC AATAAGCAACC 17520

168		
GGTTATTAAG GCTGCACAGC ATGGTAAAAA TGTTTCTGT GAAAAACCAA TTGCGCTTTC	17580	
TTATCAAGAT TGTCGCGAGA TGGTAGATGC GTGTAAAGAA AACAATGTAA CCTTATGGC	17640	
AGGACATATT ATGAATTCT TTAATGGTGT TCATCATGCA AAGAAGCTCA TTAATCAAGG	17700	
AGTTATCGGA GACGTTCTAT ATTGTCTATC AGCTCGTAAT GGTTCGGAAG AACACAACC	17760	
GTCAATATCA TGGAAAAAA TTCGTGAAAA ATCAGGTGGT CACTTGTATC ACCAATCCA	17820	
TGAATTGGAT TGCGTTCAAT TCCTTATGGG GGGCATGCCT GAAACTGTAA CCATGACAGG	17880	
TGGAAATGTG GCCCATGAAG GTGAACATTT CGGTGATGAA GATGATATGA TTTTGTCAA	17940	
TATGGAATTT TCTAATAAGC GTTTTGCCTT GTTAGAATGG GGTTCAGCTT ATCGTTGGG	18000	
TGAACATTAT GTCTTAATCC AAGGAAGCAA AGGTGCCATC CGCTTAGACT TATTCAACTG	18060	
TAAAGGAATC CTTAAGCTAG ATGGGCAAGA AAGCTATFTC TTGATTACAG AATCGCAAGA	18120	
AGAAGATGAT GATCGGACTC GTATCTATCA TAGTACAGAG ATGGATGGAG CAATTGCTTA	18180	
TGGTAAACCA GGTAAACGTA CTCCTATTATG GCTATCATCT GTCAATTGATA AAGAAATGCG	18240	
CTATCTGCAT GAGATTATGG AAGGAGCTCC AGTATCAGAA GAATTTGCAA AACTTTTGAC	18300	
AGGTGAAGCT GCCCTAGAAG CAATTGCTAC TGCAGATGCT TGTACCCAGT CTATGTTTGA	18360	
AGATCGCAAA GTAAAAATGT CAGAAATTGT AAAATAAAT TTGGTATTCT CTATTATTATA	18420	
GGTCGACTTG CTCCTCTGAA AGTACTTTTA GAGGAGCTGT TTGACTTTGC TAGTTTTGA	18480	
AACGTAAATC TATTATACTA CAAACTATTG AAAGCGTTT AATTTTAAGG TATAATAATC	18540	
TCATAGAAAT AAAGAAAAGG AGGAAAAGG ATGCCACAGA TTAGCAAGA AGCGTTGAT	18600	
GAGCAATCA AAGATGGAAT CATCGTTTCT TGTCAGGCTC TTCTCATGA ACCGCTTAT	18660	
ACAGAAAGCG GAGGGGTGAT TCCCTTGCTG GTCAAAACCG CTGAGCAAGG TGGAGCAGTC	18720	
GGTATCCGAG CAAACAGTGT TCGCGATATC AAGGAAATTA AGGAAGTCAC TAACTTCCA	18780	
ATCATTTGGGA TTATCAAAAG TGATTATCCA CCTCAGGAAC CCTTCATCAC GGTACTATG	18840	
AAAGAAGTTG ATGAATTGGC AGAAGCTGAC ATCGAGGTGA TTGCTCTGGA TTGTACCAAG	18900	
COTGAACGCT ACGATGGTTT GGAATTTCAA GAGTTCACT GTACAGGTAA GGAAGAAAT	18960	
CCTAATCAGC TTTTGTGCG TGATACTAGT ATCTTCGAG AAGGGCTAGC AGCTGTAGAA	19020	
GCAGGAATTG ACTTTGTGCG AACAACTTA TCAGGCTACA CATCCTACAG TCCAAAAGTA	19080	
GACGGTCCAG ATTTTGAATT GATTAAGAAA CTCTGTGATG TCGGTGTAGA TGTCAITGCA	19140	
GAAGGAAAAA TTCATACACC AGAACAAGCC AAACAAATCC TTGAATATGG AGTGGAGGC	19200	
ATCGTTGTTG GTGGCGCCAT TACTAGACCA AAAGAGATTA CAGAACGCTT CGTTGCTAGT	19260	
CTTAAATAAG ATGTGAGGGG GAGTTTATG TTTAAAGTTT TACAAAAAGT TGGAAAAAGT	19320	

TTTATGTTAC CTATAGCTAT ACTTCCTGCA GCAGTCTAC TTTTGGGAT TGGTGGTGCA 19380  
 CTTTCAAACC CAACCAOGAT AGCAACTTAT CCAATACTAG ACAATAGTAT TTTTCAATCA 19440  
 ATATTCCAAG TAATGAGCTC TGCAGGAGAG GTTGATTCA GTAATTTGTC ACTACTTCTC 19500  
 TGTGTGGGAT TATGATTGG CTTAGCGAAA CGAGATAAG GAACCGCTGC GTTAGCAGGA 19560  
 GTAACTGGT ACTTAGTTAT GACTGCAAG ATCAAAGCTT TGGTAAACT TTTTATGGCA 19620  
 GAAGGATCTG CAATTGATAC TGGAGTTATT GGAGCATTAG TTGTGGGAAT AGTTGCCGTA 19680  
 TATTTCGACA ACOGATATAA CAATATTCAA TTACCTTCGG CTTTAGGATT CTTTGGAGGT 19740  
 TCACGCTTCG TTCTATTGT TACATCGTTC TCTTCTATCT TGATTGGCTT TGTCTCTTTT 19800  
 GTTATTTCGC CACCTTCCA ACAACTTCTT GTTCTACAG GTGATATAT TTCTCAGCGC 19860  
 GGTCCAATTG GAACCTTCT ATATGGATT TTAATGAGAC TTCTGGAGC AGTAGGCTTA 19920  
 CATCATATAA TTTACCTTAT GTTTTGGTAT ACTGAACCTG GTGGTGTGA AACTGTTGCA 19980  
 GGAACAACAG TGGTGGAGC TCAAAAAATA TTTTGTGCTC AATTAGCGGA TTTGGCCCAT 20040  
 TCTGGAATAT TTACAGAAG AACAGGTTT TTTGCAGGTC GTTCTCAAC AATGATGTTT 20100  
 GGTTTACCGC CTGCGTGTAT AGCGATGTAC CATAGTGTTC CTAAAAATCG TCGTAAAAAA 20160  
 TACGGGGTTC TGTTTTTGG AGTTGCTTTA ACATCTTTTA TTACCGGTAT TACAGAACA 20220  
 ATTGAATTTA TGTTTCTATT CGTCAGTCCG GTTCTATATG TTGTTCAAGC ATTCTTTGAT 20280  
 GGTGTATGCT TCTTTATGTC AGACGTCTTA AATATTCAA TAGGAACAC ATTTTCAGGA 20340  
 GGTGTAAATG ATTTCACTTT ATTTGGAATT TTGCAGGGA ACGCTAAGAC GAATTGGGTT 20400  
 CTTTCAATTC CATTTGGACT TATTTGGAGT GTTTTGTATT ATATTATTT TAGATPGTTC 20460  
 ATTTACTAAT TCAACGTTCT AACGCCAGGG CGAGGAGAAG AAGTAGATT TAAAGAAAT 20520  
 TCTGAATCCG CAGATTCAAC TTCAAATACT GCAGATTATT TAAACAGGA TAGCCTACAA 20580  
 ATTTATCAGAG CCTTGGGTGG ATCAAAATAT ATAGAAGATG TAGATGCTTG TGTGACACGT 20640  
 TTACGTGTAG CTGTAAAGA AGTTAATCAA GTTGATAAG CACTTTTAAA ACAAAATGGT 20700  
 GCAGTTGATG TCTTAGAAGT GAAGGTGGC ATTCAGCAA TCTATGGAGC AAAAGCAATC 20760  
 TTAATATAAA ATAGTATTA TGAATTTTAT GGTGTAGATG ATTAAGTACT TACTGACTTA 20820  
 ATAAAAACA GAGGAGAGTG ATGGATGAGT AGGATGAAT GAAATGCAAT ACAAGAAATA 20880  
 AAGAACTCAT TATCCAAAGT GGATACGCTT ATTACATAG AGAATACAA TGAATTTAG 20940  
 AAAATPAGCT TGTACAGTAC TTGCGGGTGC TGGGTCTCTT GGTCTTGCTG CTTGTGCA 21000  
 TTCTGCCGA AGTAAGATG CTGCCAAATC AAGTGGTGAC GGTGCCAAA CAGAAATCAC 21060

170

TTGGTGGGCA TTCCCAGTAT TTACCCAAAGA AAAAAGTGGT GACGGTGTG GAACATTATGA	21120
AAAATCAATC ATCGAAGCGT TTGAAAAAGC AAACCCAGAT ATAAAAGTGA AATTGGAAAC	21180
CATCGACTTC AAGTCAGGTC CTGAAAAAAT CACAACAGCC ATCGAAGCAG GAACAGCTCC	21240
AGAAGTACTC TTTGATGCAC CAGGACGTAT CATCCAATAC GGTAAAAACG GTAATTTGGC	21300
TGAGTTGAAT GACCTCTTCA CAGATGAATT TGTTAAAGT GTCAACAATG AAAACATCGT	21360
ACAAGCAAGT AAAGCTGGAG ACAAGGCTTA TATGTATCCG ATTAGTTCTG CCCCATTTCTA	21420
CATGGCAATG AACAAAGAAA TGTTAGAAGA TGCTGGAGTA GCAAACCTTG TAAAGAAGG	21480
TTGGACAACG GATGATTTTG AAAAAGTATT GAAAGCATT AAAGACAAGG GTTACACACC	21540
AGGTTCAATG TTCAGTTCTG GTCAAGGGGG AGACCAAGGA ACACGTGCTT TTATCTCTAA	21600
CCTTTATAGC GGTTCTGTAA CAGATGAAAA AGTTAGCAAA TATACAACCTG ATGATCTCAA	21660
ATTCTGTAAA GGTCTTGAAA AAGCAACTAG CTGGATTAAA GACAATTTGA TCAATAATGG	21720
TTCAACAATT GACGGTGGGG CAGATATCCA AAATTTGGC AACGGTCAA CATCTTACAC	21780
AATCCTTTGG GCACCAAGTC AAAATGGTAT CCAAGCTAAA CTTTTAGAAG CAAGTAAGGT	21840
AGAACTGGTA GAAGTACCAT TCCATCAGA CGAAGGTAAG CCAGCTCTTG AGTACCTTGT	21900
AAACGGGTTT GCAGTATTCA ACAATAAAGA CGACAAGAAA GTCGCTGCAT CTAAAGAAAT	21960
CATCCAGATT ATCGCAGATG ACAAGGAGTG GGGACCTAAA GACGTAGTTC GTACAGGTGC	22020
TTTCCAGATC CGTACTTCAT TTGAAAAACT TTATGAAGAC AAACGCATGG AAACAATCAG	22080
CGGCTGGACT CAATACTACT CACCATACTA CAACATATT ATGGATTGTT CTGAATGAG	22140
AACACTTTGG TTCCCAATGT TGCATCTGT ATCAATGGT GACGAAAAAC CAGCAGATGC	22200
TTTGAAAGCC TTCACTGAAA AAGCGAAGCA AACAATCAAA AAAGCTATGA AACAATAGTC	22260
CTTAGTTATT CTATAAAAG TAGTTTTTA AAGAACCTAA GAGGTATAC CCCCTTTTCC	22320
CTCTACACAG ATAGTGTAAG AAAAGGGGGC TTTTGTTTAA AATGTAAAGAA ACTGTACAGA	22380
AATTAAATG AAGTTCTTAC ATAAGCGAAT CATAAAAAAT TTCATTTTGA TTTTAAACA	22440
GTTCAGAAA GTCAAAAAAT TATTCTATTT GAAAGAGAGG TGCCGACTGT GAAAGTCAAT	22500
AAATCCGTA TCGCGGAAAC AGTGATTTC TACGCTTTCC TAGCACCAGT ATTATTCTTC	22560
TTTGTCATCT TTGTGTGGC TCGATGGTG ATGGGCTTCA TTACAAGTTT CTTTAACACT	22620
TCAATGACTA AATTGAGTT TGTAGGCTTG GATAACTATA TCCGTATGTT TAAAGATCCT	22680
GTCTTTACAA AATCTCTGAT TAACACAGTT ATTTTGGTTA TTGGATCTGT ACCAGTTGTT	22740
GTCTTATCT CACTCTTGT AGCATCTCAG ACCTATCATC AAAATGTCAAT TGCCAGATCC	22800
TTCTACCGTT TCGTCTCTT CATTCTGTT GTAACGGGTA GTGTGCGGT GACAGTTGTT	22860



TGGAATGGA	TTTATGACC	ACTATCAGG	ATTCTAACT	TGTCCTTAA	GTCCAGCCAC	22920
ATCATCAGCC	AAAACATTC	TGGTTGGGA	GATAAAACT	GGCATTGAT	GGCGATTATG	22980
ATTATTCTCT	TGACCACTC	AGTTGGTCAG	CCATCATCC	TTTATATCGC	TGCCATGGGG	23040
AATATTGACA	ATTCACTGGT	TGAAGCGCG	CGTGTGATG	GTGCAACTGA	GTTTCAAGTT	23100
TTTGGGAAG	TTAAATGGCC	AAGCCTTCTT	CCAACAATC	TTTATATGTC	AATCATCACA	23160
ACAATTAAC	CAITCCAGTG	TTTCGCCCTTG	ATTCAGCTTT	TGACATCTGG	TGGTCCAAAC	23220
TACTCAACAA	GTACCTTGAT	GTACTACCTT	TACGAAAAAG	CCTTCCAATT	GACAGAATAC	23280
GGCTATGCCA	ACACAATTGG	TGTCCTTCTG	GCAGTCATGA	TGCTATCGT	AAGCTTTGTT	23340
CAATTTAAAG	TACTTGGAAA	CGACGTAGAA	TACTAAAGAA	AGGAGACAGC	TATGCAATCT	23400
ACAGAAAAAA	AACCATTAAC	AGCCTTTACT	GTTATTTCAA	CAATCATTTT	GCTCTTGTTG	23460
ACTGTGCTGT	TCACTTTTCC	ATTCCTACTGG	ATTTTGACAG	GGCATTCAA	ATCACAACCT	23520
GATACAAATG	TTATTCCTCC	TCACTGGTTC	CCTAAAATGC	CAACCATGGA	AAACTTCCAA	23580
CAACTCATGG	TGCAGAACCC	TGCCCTTGCA	TGGATGTGGA	ACTCAGTATT	TATCTCATTTG	23640
GTAAACATGT	TCTTAGTTTG	TGCAACCTCA	TCTCTAGCAG	GTTATGTATT	GGCTAAAAAA	23700
CGTTTCTATG	GTCAACGCAT	TCTATTTGCT	ATCTTTATCG	CTGCTATGCG	GCTTCCAAAA	23760
CAAGTTGTCC	TGTACCATTT	GGTACGTATC	GTCAACTTCA	TGGGAATCCA	TGATACTCTC	23820
TGGGCAGTTA	TCTTGCTTT	GATTGGATGG	CCATTCGGTG	TCTTCTCAT	GAACAGTTT	23880
AGTGAAATA	TCCCTACAGA	GTTCCTTGAA	TCAGCTAAAA	TCGACGGTTG	TGGTGAGATT	23940
CGTACCTTCT	GGAGTGTAGC	CTTCCCAGTT	GTGAAACCAG	GGTTTGACAG	CCTTGCAATC	24000
TTTACCTTCA	TCAATACTTG	GAATGACTAC	TTTATGCAAT	TGGTAATGTT	GACTTCACGT	24060
AACAAATTTGA	CCATCTCACT	TGGGGTTGCG	AOCATGCAGG	CTGAAATGGC	AACCAACTAT	24120
GGTTTGATTA	TGGCAGGAGC	TGCCCTTGCT	GCTGTTCCAA	TCGTACACAGT	CTTCTAGTGC	24180
TTCCAAAAAT	CCTTCACACA	GGTATTAAT	ATGGGAGCGG	TCAAAGGATA	ATACTCTGGG	24240
AAATCTCTT	CAAACTACGT	CAGCTTCACC	TTGCCATACT	TAAGTATGTC	CTGCGGTTAG	24300
CTTCTCATGTT	TGTTCTTCAA	TTTTCAATTGA	GTATAGGAAA	ATCAATCTAT	CAAGATACAG	24360
AAGTATATTT	TATAGATTTA	GAGAAATAG	AGGTTATATAG	TGCTACAAA	ATGGAGGGTA	24420
TGCAGTTACT	TTATGAAGTT	TTGTGACACA	CTTATTAJACT	TAAGAATGGT	TTTAGTTAAC	24480
TATCAGAAAC	GAGGAAGAAG	GTATGATTTT	TGACGATTTT	AAAAACATCA	CCTTTTACAA	24540
AGGCATTCAT	CCTAATTTAG	ACAAGGCTAT	CGACTATCTC	TACCAACATC	GTAAGGATTC	24600

TTTCGAATTA	GGAAAGTATG	AFATTTGATGG	AGATAAAGTC	TTTCTAGTTG	TTCAGGAJAA	24660
TGTCCTCAAT	CAAGCTGAAA	ATGATCAATT	TGAGTATCAT	AAGAACTATG	CAGATTGCA	24720
TTTGCTGGTA	GAAGGACATG	AATATTGAG	CTACGGTTCA	CGTATCAAAG	ACGAGGCAGT	24780
AGCATTGCA	GAGGCGATG	ACATTGGCTT	TGTTCAATGT	CATGAACACT	ACCCACTCTT	24840
GTTCGGTTAT	CACAATTTTG	CGATTTCTT	CCAGGTGAG	GCACATCAGC	CAATGGTTA	24900
TGCAGGCATG	GAAGAAAAGG	TTGAAAATA	TCTCTTTAAA	ATTTTGATTG	ATTAAAAATA	24960
GGATGAATTG	TTTTTTTGTA	AAGCTTTGAT	AATACTCTAC	CATGAATTTG	ATCTTTGTGA	25020
GGTAGAGAAA	TGAGAATAAA	ATATTTTAAA	ATTGGTATCT	TCTAAGTATG	CTGCAAGAGC	25080
TAGTTTCTTA	GATGGACAGG	GGATTACAGT	TGATGAGATG	GCTTGGATAA	TTAGGGGCAT	25140
TGTGAATGCA	TTGATTGGTA	GATACATAAA	ATTAGGTACT	TATGCGGCTA	AGTATGGTAT	25200
TAGTATGGCA	CGCTGATCT	TAAGTAGGGT	AGCTGCAACT	GCAGGCAGAA	GAGTAGGATT	25260
ACTGACCAAG	ATTTCTGAT	GGATTTTACG	AGTAGCTGTG	AAATAGGCTG	ATGATATGG	25320
TAATTTTGCC	AACAATATG	CTGCAGCTTG	GGATGCATAT	GATAAAATTC	CTAACAATGG	25380
TCGTATAAAC	TTTTAAAATG	CGAGAAATGAA	AGCACTTTGT	ATTTTTTTAT	TGAATATGTT	25440
AGCTTGGACA	GTGCTTGCAA	TGATAATTCC	TGGAGGGCTA	GATGGATTG	ATAGGCATAC	25500
TTGGAGTACT	ATTTTAATTG	CGTCGCTGTT	CGGGGTATAT	GATTATAAGC	CCATAGATAA	25560
AAATAGAJAA	AACTCCAAAA	GAAAAAATAG	ATTTGTTCAT	GGTAGGGACT	TATGAAAAGCT	25620
TTACTGACAA	AAAAGAAAAA	AGTTTACAAA	GAAAAATGAT	GGAGGAGCAA	ACATGGCACA	25680
AAAAGGAGTA	AGCTTATCA	AGGCAGCATT	TGATACAGAT	AACCTTCTCA	TGGGTTTTAG	25740
TGAGAAAGCT	TTGGACATCG	TGACAGCCAA	TCTTCTTTTT	GTGCTCTCTT	GTTTACCAT	25800
CGTGACGATT	GGAGTGGCTA	AAATCAGCCT	CTACGAGACC	ATGTTGGAAG	TTAAGAAGAG	25860
CAGACGGGTG	CTGTTTTTTA	AAATCTATCT	AGATCTTTC	AAGCAAAATC	TGAAACTAGG	25920
TCTTCAGCTG	GGTTTAATGG	AGTTAGGAAT	TGTGTTTCTT	ACCCTTTCAG	ATCTCTATCT	25980
TTTCTGGGGT	CAAAAGCTC	TGCCCTTCCA	ATTGCTGAAA	GCCATTGTTT	TAGGTATCT	26040
GATTTTTCTT	ACTATCGTGA	TGCTGGCTAG	TTACCCATAT	GCGGCACGTT	ATGACCTATC	26100
TTGGAAGGAA	ATTCCTCAAA	AAGGATIGAT	GTTCGCTAGT	TTTAACTTTC	CTTGTTCTCT	26160
CTCATGTTA	GCCATTCTTG	TCTCATTTGT	GATGGTCTCT	TATCTGTCCG	CCCTCAGTCT	26220
ACTCTTAGGT	GGCTCAGTCT	TCCTACTTTT	TGGGTTTGGG	CTATTGGTCT	TTATCCAGAC	26280
TGGATTGATG	GAGAAAAAT	TCGCAAAATA	CCAAATAGGAG	CTTTATTTCT	GAAACTACTT	26340
TCAAAGGCTC	CAAAGCTAT	TCTATAAGCG	AGAAACTAAA	ATCCG		26385

## (2) INFORMATION FOR SEQ ID NO: 4:

## (i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 2716 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 4:

CCTGCCCGCA TTCCCTAGG CATTAACTAA ACATATAAAA GCATGTGAGA GACTGTTGGA	60
AAAGCGAGGA AATTCCCTCT CTTTCCTCT AGTCTCTCT TCTTTTGTCT GATTTTATTC	120
AAAGAAATG ATATAATAGT AGTTATGAG AAAAAAGAAAT TACGCATCAA TATOTTGAGT	180
TCAAGTGAGA AAGTAGCAGG ACAGGGAGTT TCAGGTGCTT ACCGTGAATT AGTTCGTCTT	240
CTTCACCGTG CTGCCAAGGA CCAATTGATT GTTACAGAAA ATCTTCCAAAT CGAGGAGAT	300
GTGACTCACT TTCATACGAT TGATTTTCCC TATTATTAT CAACCTTCCA AAAGAAACGC	360
TCAGGGAGAA AGATTGGCTA TGTGCATTTT TTGCCAGCTA CACTTGAAGG AAGTTGAAA	420
ATTCCATTTT TCTTAAAGGG AATTGTGAAA CGCTATGAT TTTCTTTTAA CAACCGGATG	480
GAGCACTTGG TTGTGGTCAA TCCTATGTTT AITGAGGATT TGGTAGCAGC TGGTATTCGA	540
CGTGAAGGAG TGACCTATAT TCCTAACCTT GTCAACAAGG AAAAATGACA TCCTCTACCA	600
CAAGAAGAGG TAGTCAGACT GCGCACAGAT CTGGGTCTTA GTGACAAACA GTTTATCGTA	660
GTAGGTGCTG GGCAGTTCA GAAACGTAAA GGGATTGATG ACTTTATCCG TCTGGCTGAG	720
GAATTGCCCTC AGATTACCTT TATCTGGCT GTGGCTCTC CTTTGGTGG TATGACAGAT	780
GGTTATGAAC ACTATAAGAA AATTATGGAA AATCCCCCTA AAAATTGAT TTTTCCAGGC	840
ATTGTATGCG CAGAGCGGAT GCGCGAATTG TATGCTCTAG CGGATCTTTT CTTGTGCCCT	900
AGTTACAATG AGCTCTTCC TATGACTATT TTAGAAGCTG CGAGTTGTGA GGCCTCTATT	960
ATOTTGCCGTG ATTTAGATCT CTATAAGTG ATTTTGAGG GAAATATATG GCGGACAGCG	1020
GOTAGAGAAG AGATGAAGA GGCTATTTTG GAATATCAG CAATCTCGCT TGTCTAAAA	1080
GATCTCAAAG AAAAGCTAA GAATATTCC AGAGATATT CTGAAGAGCA TCTGTTACAA	1140
ATCTGTTTGG ACTTTTATGA GAAACAAGCC GCTTATGAGA GAAAGTAAA AGTGAGGTAA	1200
TCTATCGGAA TTGTTTTATT TACAGATACC TATTTTCTCT AGGTTTCTGG TGTGCGACC	1260
AGTATTCGAA CCTTGAAAAA AGAAGTTGAA AAGCAGGAGC ATGCTGTTTT TATCTTTACG	1320
ACGACAGATA AGGATGTCAA TCGCTACGAA GATTGGCAAA TTATCCCATC TCCAAGTGT	1380

174

CCCTTCCTTG	CTTTTAAGGA	TCGTCGCTTT	GCCTACCGAG	GTTTTAGCAA	GGCACTTGAA	1440
ATTGCTTAAC	AGTATCAGCT	AGATATATATC	CATACCTCAGA	CAGAAITTTTC	TCTTGCCCTG	1500
TGGGGGATT	GGATTGCGCG	TGAAATTGAAA	ATTCCAGTCA	TCCATACCTA	TCACACCCAG	1560
TATGAAGACT	ATGTCCATTA	TATTGCTAAG	GGGATGTGA	TCCGGCCGAG	TATGGTCAAG	1620
TATCTGGTTA	GAGGTTTCT	GCATGATGTG	GATGGGGTTA	TTTGCCCTAG	TGAGATGTGTC	1680
CGTGACTTGC	TATCTGATTA	TAAAGTCAAG	GTGAAAAAC	GGGTCAITTC	TACTGGGATT	1740
GAATTAGCCA	AGTTTGAGCG	TCCGAAAATC	AAGCAGGAAA	ATTTGAAAGA	ACTGCGTAGT	1800
AAACTAGGGA	TTCAAGATGG	TGAAAAGACG	TTGCTTAGTC	TTTCGAGAAT	CTCCTATGAA	1860
AAAAATATTC	AAGCAGTTT	AGCAGCCTTT	GCTGATGTTG	TGAAAGAGGA	AGACAAGGTT	1920
AAACTGGTAG	TAGCTGGGGA	TGGCCCTTAT	CTGAATGACC	TCAAAGAGCA	AGCCCAGAAC	1980
CTAGAGATT	AAGACTCAGT	CATCTTTACA	GGGATGATTG	CTCCTAGTGA	GACGGCTCTT	2040
TACTATAAAG	CGGCGGATT	CTTCATTTTCG	GCATCGACAA	GCGAAAACGA	AGGTTTGACC	2100
TACTTGGAAA	GCTTAGCCAG	TGGAACACCT	GTCAATGCTC	ACGGAATTC	TTATTTGAAC	2160
AACCTCATCA	GTGATAAAAT	GTTTGGAAAC	TTGTACTATG	GAGAACATGA	TTTGGCTGGT	2220
GCTATTTTGG	AAGCCCTGAT	TGCAACACCA	GACATGAACG	AGCATACCTT	ATCAGAGAAA	2280
TTGTATGAGA	TTTCAGCTGA	GAACCTTGGG	AAACGAGTGC	ATGAGTTTTA	TCTGGATGCC	2340
ATTATTTCAA	ATAACTTCCA	GAAAGATTTG	GCTAAAGATG	ATACGGTCAG	TCAGCGTATC	2400
TTTAAGACAG	TTTTGTATCT	TCAGCAACAG	GTGGTTGCTG	TACCTGTAAA	AGGATCTAGA	2460
CGCATGTTGA	AGCCTTCAAA	AACACAGTTG	ATCAGTATGA	GAGACTATTG	GAAAGACCAT	2520
GAAGAATAGA	AAGAGGAACA	GCTATGAAAA	AAACAATTAA	TGAGAAAGCG	TCGTGATAAA	2580
AAGATTGCGG	GTGTTTGTGC	TGGGGTGGCC	CATTATCTGG	ATATGGATCC	GACTATCGTT	2640
CAAGTCATTT	GGGGTGTCT	TACTTGCTGT	TACGAGCTG	GAATTGTAGC	TTACATTATT	2700
TTATGGATTA	TCGCGA					2716

(2) INFORMATION FOR SEQ ID NO: 5:

- (i) SEQUENCE CHARACTERISTICS:  
 (A) LENGTH: 13926 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 5:

CTTTGGTTT GCCTATTCA AGACATGAGG GCCATCAGGA ATGATCTGAA ACTGCGAATC

60

TGTTAACAGT	CTATGGAGAG	CTTTCATAGA	ACTAAGATTC	GGTTTATCTT	TCCTGCCACA	120
AATTAGTAAG	GTTGGATAAG	GGTAAGTTCC	TGCTATATCC	GTTAAATCAA	GTGTCTTCAA	180
CTCCTCAGAA	ACTCCGACCA	TAAGAGTCTT	GTCTGTCTCC	TGTTTTTCAA	ATACTCTTTT	240
GGGAAGTAGT	TTAAAATCA	GCAATTGAAG	ATAAAATAGG	ATATTTCCCTG	CTAATTTAAG	300
CGGGATCTCT	GACAGATCA	AAGCTCGAAG	ATTGTGTAAA	TCGTAACTGG	AAAGTTCAG	360
TGTCAGGGCA	GCACCTAAGG	ACAATCCAAT	CAAAACAAAA	GGTTCGTCT	CTTGAGCTAG	420
GTCTGTGATA	ACTCGTCTTT	TAGCTTGTGT	ATAGTTACTA	ACTCCAGAAG	GAAATAACTC	480
GATAGCCTCA	GAAGGATAAT	CTGTTCAGTAG	ATTCCGAAC	TCTTTCCAAAG	ACTCTGCTGA	540
CTGCCCTAAC	CCATGCAAAA	ATATTTAATTT	CATCTAGTTC	TCCTCAAGGC	TTAATTCATA	600
CAAGCCTCTC	ACTGCATTAC	AGCCGTAAAT	AGCTTCTGCT	TGGGTAAAT	CTGCCAAGGT	660
CAAGACTTTC	TCTTCTACCT	GTCTGTCTTC	TAGCAATGC	TGACGGTAAA	TTCCTGGCAA	720
GATTCCAAGT	CGGATAGCGG	GTGTGTAGAG	TTTCCAGCG	ATTTTCAGAA	CCAAATTTCC	780
TATAGAGGTT	TCAAGCAGTT	CTCTTGACTT	ATTGTGGTAA	ATCTTCTCTT	GTTCCTCTAG	840
GCTCAAAATGC	GGTCGGTAG	TGGTTTTAAA	GTAGGTAAAG	GATTTGATCA	AAGCAGCTTC	900
GTGAAGACAG	ACTTGGGCT	GACAAAAGCT	TGTACTGAGA	GGGGTTAAAT	CTTGACGATT	960
GACTTCTATC	TCTCCAGATT	TGCTAAGGCT	GATTCCGAA	CGGTAACTTC	GATTAGCTTC	1020
ACAATCTCGA	CACCTTCTCT	CAATCTTGTG	TCCCAAGTCT	TCTGCATCAA	AAGGAAAAGC	1080
AAAATPACGA	CTAGCTTTTC	TCAGCCTTTC	CAGATGTTGT	TCTTCAAACA	TCAGTTGTTT	1140
TTGGCTGATT	TTTCCAGTTG	TAATTAATTG	GAAGCGAGCT	TGTTTACGAT	AGAGAACTGC	1200
TGCTTTTGA	TGAACCTCTC	GGTATTCAGA	TTCCCATGTG	CTATCCCAAG	TAATCCCTCC	1260
GCCAACTCCA	TAAATGGCTT	GACCTTTGTG	AAGTTGAATG	GTACGAATGG	CCACATTAAT	1320
AATCCGTCTG	CCATTTGGAA	GCAAGAGACC	AATCGTTCCA	CAGTAGACTC	CACCGGGTTG	1380
AGGCTCCAA	TCTTGTGATA	TCTCCATTGT	CGCAATTTTC	GGTTCACCCG	TTATGGAAAC	1440
ACAAGGAAAG	AGTGAGCGGA	AGATTTCAAC	AAGGTCCACA	TCTCTCGCA	ACTGACTCTT	1500
GATGGTCGAA	GTCTCTGCCC	AAACAGTTGA	ATACTGCTCT	ACCTGACACA	GACGCTCCAC	1560
GTCTCTGCTC	CCAACCTCAG	AAATACGGTT	CATATCATTTG	CGCAAGAGGT	CCACATTCAT	1620
CATATTTTCA	GAGCGATTTT	TGGGATCCTG	TTCCAAACCA	CTGCCCTGTT	CAAGATCTTC	1680
TTGGTCAGTT	ACCCCAAGCT	GAGTCGTCCC	CTTCATTGTT	CGTGTGTGTA	ACTCGGATC	1740
ATTTGTCTCA	AAAAGAGCT	CTGGCTCAT	GGAAATCACT	GTCACTCTGT	CATGTTCCAC	1800

176

ATAGGCATTG TAGCCCGCCT CCTGCTCTAC CACCATACGA TTGTAGATGG CAAAGGATT	1860
GGCAITTTAAC TTTTGCTTAA GTTGGACGGT GTAGTTGACC TGATAGGTAT CTCCCTGCCG	1920
TAAATGATGG TGAATTTGGG CAATGGCCCT TTCAATAGTCT GCTGCAGACG TTACTTCCTG	1980
CCAATTTGAG GGCAATCAA TATCCCTATA AGTCAGAGGA ATAGGGGAAG TTCTTACGAT	2040
ATCATGAACTA GTAAAGTAAA GCAGGTACTC TCCCACTAGG GGATCCTGT GAACCTGCTAA	2100
TTTTTCTCCA AAAGCAGGTG CAGCCTCGTA GCTGACATAC CCCACCACAT AATAACCTTG	2160
CTCTTGCTAG CTTTCCACTT GTGCCAGCA ATCTGCCACT TCTTCTACAT TTCTCGTTTT	2220
CAACTCTTTA ATAGGCTGGG TAAAGGTATA TCTCTCCGCC AAAGTCTTAA AATCAATCAC	2280
TGTTTTCTTA TGCATACCTT AAGTATAGCA TAAATAAGA AAACCTCAT CCGCAAGCA	2340
GATGAGAGAT TCAATTTATT TAAAGATTGA AGTTTTAAG CTATTTGTTT GTTGAAGAAG	2400
TTCTTATATA ACAGCTTCTT TTAATTTAAC TGTATATTTC ATAGATACTG TTTTATTACC	2460
GTTTCTCTGT TGTTTAAGAG TTTCCGCATC TTTTTTAACA GCTTCTTTAA ACAATGTCAG	2520
TAAATCATCG TATGATGAAA CGGAAGAACC ATTACTTCG AATGTTGTTA ATCTTTTCGT	2580
TGCTTTATCT TTAACTTCTT TGAAGTAAGC TTTTTTAAAT TCTTCAATAG TATTAATGT	2640
ATTGTTAGAT ATTTTCTTGA TAATATATTC ATCACTTAGA ACAGACTCAC CATCTGTTTT	2700
AGATTGTTGT TTATATTTAT TTGAAGCATA ACCTAAGAAC CCATTTCGT ATCCGTAGTA	2760
ACCCCAATAT CTAAGAAGAT TATGTTTAAA TGAACAGCT CCAGGAGCAC CTTTACTAGT	2820
ATTACCTCCG TAGATACCGG TCATCATTTCT AACACCTACA TAAGGTGATT GATCGTTATA	2880
GCTAATTCGT TCGGGTTAT AGATACCAAT ACCTGGATTG CGATTAGTCA TTAATGTTG	2940
ATCAACTAAA TCATTACAG ATTGAAATTT TAATTCATTT TTCTCTCTCT GACTTAGATT	3000
TCGAATTTTA TCCCATGAT TTAATTTATG GTTATCAAGG TATTCCTAT CTATTTTTTT	3060
GAACCATGCA CTATTTAAAT CTTTATTTTG TTGAGAAATC ACAGATTCAG CCTCAATTC	3120
ATCAAGAAGA GTTAAAGTGT CATTAATAAC CTTCATATAT CTATTAATAT CTTCTCGTGT	3180
TTTTAGAGTT TTTGGATCTG TAATATAACA CTGATTCOCA TCAITTTTGC GTTTAAATAC	3240
CATATTAAATA CCTAAGAAC CAAACTCATC AAATCCACTA CCAATACAG GAGTTGTAG	3300
CATACCTCGA GCATATGCTT CAGCATCAGT ACCTTCACGG TGTCCAAGC CACCTAAGTA	3360
AATCGCACGG TCGTTGACGT GTGTGTTTTC ATGTGTGTAA ACTGAATAC GTATTCCACC	3420
AACCATTTCT AAATGAACAT ATTTTACATC AGTTCTAATA TCATCAGAGT TAGGATATAT	3480
AGCAGCATAA GCTCTGTTT CATTAATAAT ATAACTACTA TCCATAGGAC CAAGAAATTC	3540
TCTAAGAGGA GTATATACTT TGTGCGTATT ATAGCGGCA TATTTTTCAA CCCATCCACC	3600

AGGAGCGTTA	TAACCTTCCC	AAATAGGAAT	AACAGCATCT	CTTAGTAGTC	GTCTTTAAAC	3660
CTTATCAGAC	GCTAGACGAT	ACCAGAAATC	ATAATAGTTT	CTATAACCAT	CTGCAGCTTT	3720
GTTAACGATA	TCTTTAATAT	CTTCTAARGA	TPTTTTACCT	AATCGCTCTG	CACTACCAAA	3780
CCCAATTGCA	TTATAATTTG	AAATTAAATA	AAGATGTGCT	TTATCAATAT	TCAGTAGTGG	3840
GAGTATAGTA	TTTCTAAGGT	GACTTCGTIT	TAAATTATCG	AATGCACGAT	GTTTAGAAT	3900
TTTAATTTCT	TCGACCTCAG	AAGCGCGTTC	TGCGATGTAG	ACATCGTCTT	CTGTAGCATT	3960
AATAAACCAA	TCGTTCTATAT	TGCTATATTT	TGIGAACAAT	TGCTCTATTAT	AATTTAAAAA	4020
TGCATCTAAA	TTACCTGAT	TAGTATATTT	AGCCAATACT	TGACCGAATG	CGTGAATGT	4080
ACGTGAACCT	TTAATGTTGT	TCTCTTTAGA	ACCGATTCCA	ATTAAATCTGT	CTAATACGCT	4140
AACTTTTTCA	CCATAGAAAT	CTGGTTTGAA	TAGCATTAAAT	TCTTTAATAT	TAACATCACC	4200
AAATTAACT	CCATAGTAAC	GATTTAGGTA	AGTTAAACCT	AGTAATAAAG	CTGCTTTGTT	4260
TTTCTCGACT	TTATCACGAA	TCATTTGACG	AGCAGCTGGA	GAATCATTTA	GTTGATGTT	4320
TTGTTTTTGA	ACTAATTTTG	TGATTAGGTT	TGTTAAGTTT	TCTTTAACAT	CTGTGAAGCT	4380
TTCTCTPAAA	TATAAATCTT	TGATTGCAT	AACCTATFAG	TCACCTAATC	GATTTAGATG	4440
CTGATACATC	GTTTGAGACT	GAAGCTCTAC	TGATTTCTAAA	ATAGATTTTA	TATCATTTAC	4500
AAGAGTAGTG	TTATCTTTTT	GAACGATATT	AGGTGATATAT	TTAAATCTTA	AGTCAGTTAT	4560
AGTATATCTT	TTTACATTAC	TTAAACCTTC	ACTGCTAGAA	GACAAGTTAA	AGTAATCTTT	4620
TGTACCGTCC	GCATAGTGAA	CAATAATTTT	ATTAGCTTCA	TCTAGGTTTG	TGATAAACTC	4680
ATTGTTGTTT	ATCGCGGTAA	CAGAAAGAAC	TTCTTTAGTA	TTTATAGTGT	GTTCITTTAT	4740
TAATTTATTA	CCTTGATATA	CAATAAATC	TTTATTTAGT	AATGGTATTA	ATTTTTCAAG	4800
ATTTTTATAG	GCTTGGTTAT	ATTCAGCGTT	ATAATCTTGA	ATACAGAAAT	AGGCTTTTTC	4860
TTCATTAAGT	TTTGCAAGAG	GAGATAGATC	ACTTTCTAAT	TTATCAGCAG	TAATATTGAA	4920
AGTATTAAGT	TTAGCATCAG	CTTGTCTTTT	AGTAATTTTA	GTAATGTTT	TAGATTTTCT	4980
AAATGATCTA	TTACCTGACG	AATATCCCTC	TACGCGATAT	AAATCTTTTA	TATGAGCACT	5040
AGCATATATCA	GAATCATCAA	CGTGGTTAGA	GCCGAATAAC	TCCTCTCCAC	GGATAATCTT	5100
AGCATAGCTG	ACAGAATTAC	TTACCGTACC	TACAGGCCAA	GTCTTACTTG	CTATTGCTCC	5160
AACCTTCTACT	GGATTTTGAAA	CATCTATTTT	ACCTTTTACA	ACCGACTCAG	TTAGGAGAGC	5220
TTTTGTGACCA	ATAAGATGGT	CTAGAGTTAA	TCCATAATCT	ACTTTTAGGAA	CTAACAAGCT	5280
GCGCGTGTTT	TTGTTTCTCTG	TAAATAGTAG	ATCAACATAT	GCTTTTCTAA	CAATTCCTCT	5340

ATAGTTTGTA	CCTGCAATTC	CCCCTGTATG	AGAGCCATTT	CCACTGTGAG	AGTGTAGTTC	5400
GCCAAAGAAA	GCAACATTTT	CAATACGAGT	TCATCATTC	ATATTATTTA	CAAATCCAGC	5460
AACATTATTA	CGACCTGAAA	GTGTGCCCTGT	AATTTTGACA	TTTGTAATAA	CTGAAGAACC	5520
TTTCATAGTA	TTGGCTAATG	ATGCAATATT	ATCTTGACCA	GAACGTTCTA	TCTCTACATT	5580
TTCAAAMTTC	ACATTTATTA	TCGTTGCCGT	TGTTATCACA	TTAAATAATG	GATGTTCCAA	5640
TTCAGTAATA	GCAAAATGTT	TTCCCTTCAGA	ACTTAAAAAGT	TTTCCGTGTA	ATTCTTTAGT	5700
GATATATGAT	TTTCCATTAG	GAACAACATT	TCTAGCGCTC	ATTGATTGTC	CCAGACGATA	5760
TTCTTTTGAA	GGATCGTTTT	GAATAGCTTC	CACATAATCT	TTGAAATTAT	AATATACATT	5820
ATCTTCGTGG	ACTTTAGGTT	TTTCAATATA	GTGAACGTAT	TCTTCTTCAA	ATTATATTATC	5880
AGCAGTTCTA	GAGACTAAAT	TGCTCGCAT	TGCTGTAAC	TTATATACAG	GTGTTCCGTT	5940
AACCGTAGTT	TCTTCTATAT	TTTTAAACAGC	TAGTANTGTA	GTTTTCTGAT	TATTTGAAGT	6000
TATTTTAAAA	TAATAATTGC	TCTTATCATC	AGGAATAGTT	GTTATCAGTG	ATTCAATTAGT	6060
TTCTTTTCCA	TTTTCGTATT	TGATTTAAATC	TGTACGTTTA	ATATTTTAA	GCTCAACTTT	6120
TTTAAGATCT	AATTGAATAT	TTTGATTTTC	TAGAGTTTCA	GTTCCTTCAC	CGTTACCTCT	6180
CTCGTAAATC	ATAGTTGTAG	ATAGGGTGTA	TTCTTTGTAG	TACTCTAGGT	TCTTAAATGC	6240
AGCGCTTATA	GTTCCTGTGT	TTACCTTGTC	ATCTGTAAGG	ACTACAGTAT	TAATAACTTC	6300
TTCTCCTTTT	TTCAATTCAG	CTGTGATTGA	TTTGATTTTC	GTTTTGTTC	GATTTTCTAG	6360
AGTATACCTA	GCAACAGCTT	CACGTTCCAA	TATTTTCTTA	TCGGTACTAG	TCAATGTTAA	6420
TATTGGCTTT	TCAGATAATT	CAACCAATTT	TTCAATAGTT	GCAGTTAATT	TTTCAACAGC	6480
TTCTGTTTCA	TCACTTTGTT	TAGCATCTGT	ATTAGCTGCA	ACTTTTTCAG	CCTTTGTAAAC	6540
TTCAGTTTGG	AGGTTTGGCC	AACTTCTATC	ACTGTAATGT	TCTTTTACCT	TTGTTTGTGC	6600
ATCTGCAMTC	GTATTGTTTA	ATTACGTTTT	ATCAACGTTT	AGAGCGTCAA	TAGCCGTTTT	6660
AAGTTTATTT	GTCTCGCTAT	TTACCTCAGG	CTGTTTACAA	GGCTCTGAAG	CATAGACACC	6720
TTTTGCAGTT	TCTAAAACAG	GTCCAAGAGC	ATTGTAACCT	GCTGTAGAAT	AATCAGTAGG	6780
AGAAACTGAA	CTAGCTTTAT	CAATTTGATT	ATTTAACTCA	CTTTATCAA	CTGGTCTCTT	6840
AGTACCAMTA	CCCTTTATTT	TATCTTCTGG	TTTCGGTGTG	TCCTCTACAG	CCTTCTCTTC	6900
TTCCAGGAAC	TCGTGTTGCT	TTTCTGGCTC	AACTGGTGCC	GTGGTGCTCT	GTTCTGCTTC	6960
TCPTGGCGCG	ACTGGTTTCA	CTGCTTGTTC	AACTTTGGGT	TCCTCTGTGT	GTTCGTGTTG	7020
TTTCTCTACA	GCAGGCGTTT	CAACTTTTGG	TTGTCTAATA	GATTGATTAA	CAGTCTCCTC	7080
TTTTGGTTCT	ACAGTTTCTT	CAGCCTTGGT	ATCTGGAGTT	GACTCTTCTT	GTTCGGGTGT	7140



TTCTCTACAG	GCTTCTCTT	CTTCAGGAGC	TTCTGGTTGC	TTTTCTGGCT	CGACTGGTGC	7200
CTTTCTGTT	TCTCTTGGCG	CGACTGGTTC	ACCTGCTGT	TCAACTTTTG	ATTCTCTCAGC	7260
TGGTTTGTCT	GATGGTTGAC	TTTCTGGCTT	AACCTGCTACT	TTTTCTCTCG	GTTTTGACTC	7320
AACCTCTCCA	CCTACTACTT	CBACTGGAGC	TGGTTCTGCT	GAATCTTCTT	TCCCTCTCTC	7380
TACTTTAGGA	AGGGTGTCT	CAGTAGGTTT	TACCTTCOGAT	TTTGGTTCTT	CTTTTGAGCT	7440
TTCTCTGTCT	TTAGGTGCTT	CTCTTTTGG	AGCTTCTCT	GTCCTACTA	CTTGGTTTTC	7500
TGCTCTAGCT	TGCTCTGAT	TTGTTATTGA	TTGAGGAGTC	TCAACTTCGA	CCACAGTCAC	7560
CTCTCCAGGT	TTTGCTGAGG	TTTCTTCTAA	AACAGTGTCC	AAGCCAAGCG	TTTTGAGGAT	7620
GTCACTTGAT	AGATAACCAA	CATAGCGATA	GCCCTCCATT	TCAACAACAC	CCTCTCGACT	7680
AGCCAGCGCT	AGGGTCTGAA	CTGGGTCTAC	AGCCCTTGCA	CTAGGAAGAA	CTACCAATCC	7740
CATAGCTCCA	ACTAGAAAGA	CGCTAGCAAT	TTTCTTTCTC	TTGTAGATTA	AAAGCAAGCT	7800
CCCAACAGTC	AGCAAAACAA	AAGCTGTCAA	AACAGATGCT	TCTGTCCCTG	TTTGAGGCAA	7860
CTGATCTTTT	TGATACACCA	AACCATATAC	AACCTTCATT	CTGTCAAGGT	TTCCGTCTCT	7920
AATTAATCT	TTAGCTTCTT	GTGAAATAAT	CTCTTTATTT	ACATAGTGAT	AGGTGGCTGC	7980
GTCCACTACA	GAAGGAGCCA	TCAAAAGGCT	TCCAAGAAAT	ACAGAGCCTA	CAACTCCCTT	8040
AATCTTACGA	ATTGAAAAAC	GGTCTTTT	AAACACTTTT	ATCTCCTTTA	TTCAATCTCA	8100
AAACTTCTTA	ATAGCATCTT	GCGGATAOTG	CGCACGCGCA	CCTCCGATTA	ATTTTGAGCG	8160
ACTAGCCAGT	GCCGTTACAT	GGGCATGACC	AATCTCTCTC	AAAATAGGGC	GAATCGGAAC	8220
CTGAACATGC	TTGACATGCA	TGCCAATTGC	AGTGTCTCCG	ATATCCAATC	CAGCATGAGC	8280
CTTGATAAAT	TCAACCTCAA	CTGGATCTCT	CATAAACTTA	AAGGCTGCCA	ACTGCCCCGA	8340
ACCTCCTGCA	TGAAGAGTAG	GATGGACACT	GACAAATTCC	AGACCAAACT	GCTGTGCCAC	8400
CTGAGCTTCA	ACAACGAGAG	CCCGATTGAC	ATGCTCACA	CCTTGAAGTG	CTAAATGGAT	8460
AGCTCTACTA	CCTAGAAAT	CCAAGATAGT	CTCCACTATC	AGCTCACCAC	TCTCTTGACT	8520
GGATCTTTT	CCAATATGAC	CACCTAGCAC	CTCACTAGAA	GATAGACCTA	AAACAAAAG	8580
GGCCCCCTGC	TTCAAAATGG	TCTTTTCTAA	AACATCTTCC	ACTACCTGAC	GTGTTTCTCT	8640
TTGAATCTGT	GTCTCGTTCA	TCTCTGTTAC	CTCTGTTGTC	ACTCTTCTAT	CATACCGTTT	8700
TTTCTTGTTT	TTAGCAAGAT	AGACAACCTA	GAAAGTTTGC	CCAATTACGC	ATAAACTCTC	8760
CAGAATTGAC	TGGGAGTTAG	CTAGTTTCTA	TTCTATTAT	ATATATTCTA	ACTTTCTGTC	8820
CTTTTGGGG	TCTAGAAATCA	ATCTTCATAT	GATAAATTGC	TCCAAAATGA	AGTTTGAGGC	8880

180

GTTGATCGAC	ATTTTGAAAG	CCAACTCCCC	CACGTTTGAG	TTGACTTTGA	CTACTATCAC	8940
CAGCATCTTG	GAGGCCAAG	CCATCATCCT	CAATACGGAT	GACCAATCCC	GAATCCCTGT	9000
TCTGGACAGA	AAGTTTAATA	TGGCCCTGAC	CTTCCTTTTC	CTTAATGCCA	TGGTAAAGAG	9060
CATTTTCTAC	AAGGGGTGT	AGGACCAAGT	TGGTAAGAC	TAAATTATCA	AAGGCAACAT	9120
TTTCATTAAT	TTGCTATTC	AGCTTATCTC	CATAGCGTTG	TTTCTGGATA	ARGAGATACC	9180
GGCGGACATG	ATTGATTTCG	TCAGAGAGAC	AAATCAAGTC	CTTGCCCTGA	TTGAGCGCCA	9240
AGCGGAATA	GGTGCCAAG	GACTTGGTCA	CCTGCACCAC	TCGCTGACTA	TCATGAAATT	9300
CAGCCATCCA	GATGATGGTG	TCCAAAGTGT	TATAGAGGAA	ATGTGGATTA	ATCTGGCTCG	9360
AAAGGGCTTG	AAGTTGGTAC	TGACGGGTG	TTTCTTCCTG	GCTACGAATA	GCTACCATCA	9420
ACTGATCAAT	CTGATCCAAC	ATAGCATTA	ATTGGCGAGT	TACTTCTCTC	AGTTCATAGG	9480
CACCAACTTC	CTTGGACAGA	AGATTTTGAG	CACCAGAAGC	AATTTCCAAC	ATGGTTTCTC	9540
TCAAATCCTT	CAAGGAGCA	ATCCAGCGTT	TAAGACTGAA	CACACATAAG	CAGAGACAGA	9600
CAAGAGAGAG	TGTGACACTG	GCCCCAAGCA	AGGTCCACAA	GAGCTGACTC	CGAACCTGCT	9660
CTAATCTTTC	CAATGATGAC	ACGCCAAGCA	CGTCCAACT	AGTTCCTGCA	ATCTTCTACT	9720
GACTGACGTA	GGATTTTGGA	CCAGGAGTAT	AACCTGACC	TGTATCGATG	TAGGGTTTCA	9780
TAGCTCCCAT	TTTGCTAGAC	GAACTATAAA	CTGTGTGTTG	AGGATCGTAG	ACAAATTCAT	9840
GGTTTTCATT	GATTAATGAAG	GCAAGGCCCT	GCTGCCCCAA	CTGGAGTTGA	TTGAGATAGG	9900
CTTCCAGAGT	TTCTATAAGAA	ATATCCAAAC	GAAGCACACC	AAGATTGGCT	CCCTTTGCAT	9960
CAACAAGTTC	TTGAGTGACA	GAAATGACCC	ACTGACTATC	TGATTTAGCA	GCTGGAGTCA	10020
AAACAGGCAT	AGCTCCCTGA	TGAATGGCCT	TTTGGTACCA	ATCTCAGCC	ATCATATCAG	10080
AGGAAGTTT	CATCTGCACA	CTGTCACTG	TAGAAATGAC	CTGACCAGAT	TTGGTACCA	10140
GCACAACAGT	TTTCAAGTCC	TTATCTGACT	TCAAGATGGT	CAAAAACAA	TCTCGGATTC	10200
CCTCGACCTT	GTCTTGACTG	GGATTCTCAG	CATAGGCCAG	AACATCCGTC	TGCTGGGTCA	10260
AACCACTCGA	GGTGTGTTCT	AGTTTCTTGA	TATAGACTG	AATAAAGTGG	CTAGTCTGGC	10320
TGATGGTCTG	TTGGCTGTG	CCCTCAATGG	TGGCTCAAT	GGCTGAAGAA	CTTGATTGAT	10380
AGTAGAAAGT	TCCAACCAGA	GCTAGGAGAA	TGAGAAAGAC	CAGAAAGATG	GAAATACCA	10440
TTCTAACTAA	AAGAGAAAGAA	CGCTTCATCG	GTCTTCTCCC	TTCTTAAACT	GACGAGGTGT	10500
CACACCTGCA	ATCTGCTTAA	AACGTTGGGT	AAAAAGTTC	ATATCTTCAA	AACCAACCTT	10560
CTCTGCGATC	TCATAAACTC	TCAGATCTGT	AGTTAAAGAC	AAGAGCTTGG	CTTGTTTAA	10620
ACGTTCTCTC	ACCAGATAAT	CCTGAAAAGG	CAAGCCCAAC	TCTTCTTAA	TCAAGGAAT	10680

CAGATAGGTC	GGACTAAAC	CTAAGTCAC	GGCTAAAGAC	TTTAACTAA	ATPGGCTATC	10740
AGCCAGATGA	GACTGGATTT	TCTGGGCCAT	GTTTCCTTCA	AACCTATTAG	TCAATAAAATC	10800
TTGTAACGTC	TCCTCTTCT	CTTCCTTGTC	TAGTCTTTGT	TGTATTTTCC	CCAACATTTC	10860
CTCAATATCC	TGACGAGAA	AGGGTTTGAG	CAGGTAGTCG	TCCACACTTA	GTTTGACAGC	10920
AGACAAGGCA	TAAATCAAAAT	CATGTAAAC	TGTTAAAAAG	ACCAAATGAA	CCTGAGGATA	10980
GGTTCTCTGT	ACCAGACTGG	CCAACTGGAT	GCCATTTAGA	TGAGGCATGT	TGATATCGGT	11040
TAAATGATA	TCTGGCACCT	GCTTTTGGAT	CAATTCCCAA	GCCTGCCTTC	CATTTTCAGC	11100
CTGACCGATG	ATTTCATAT	CGTAGGCTGC	TACATTGACC	AGTTTAGTCA	AACCTTGCTT	11160
TACCAGATAT	TCACTCTCTA	CGATTAAAGT	TGTGTAGGTC	ATGCTCTGCT	CCTTTACCAC	11220
TTACTAGTAT	CAGTAGAGCA	AAATTCCTCT	CTAACTGCTT	AGGAAAGACC	TCTTATACTC	11280
AATAAAATC	AAAAAGTAAA	CTAGGAAGAT	AGCCACAGGT	TTCTCAAAGT	ACCGCTTTGA	11340
GGTTGTAAAT	AAACTGACG	AAGTCGACTC	AAAGTATAGC	TTTGAGGTTG	TAGATAAAAC	11400
TGACGAAGTC	GATAACCTTA	CATACGGTAA	GGCGACGCTC	ACGTGGTTTG	AAGAGATTTT	11460
CGAAGATAT	TAATCAACAT	AATCTAGTAA	ATAAGGCTAC	CTTTTCTTTC	CATTTGGTCT	11520
TTGGGAATAA	AGCGGATAGA	GAGGCTATTG	ATACAGTAA	GTAAGCCGCC	CTTGTCTCTG	11580
GGACCATCCG	TAAAGACATG	CCCAAGGTGA	GAATCTCCTA	CTCGGCTCCG	CACCTCCATA	11640
CGGTCATAT	TGTAGGACTT	ATCTTCCCTG	TAGGTGACAA	CATCTGGACT	GATGGGTTGG	11700
GTAATACTAG	GCCAGCCACA	ACCAGACTCA	AATTTGTCCT	TTGATGAAAA	GAGAGGTTCC	11760
CCAGTTGCTA	TATCCACATA	GATACCGGAT	TCAAATTTAT	CCAGTAAACG	GTTTGAGAAA	11820
GCTCGTTCTG	TTTGATTTTC	CTGGGTAACT	GCATACTCCT	CAGGTGACAG	GGTCTTTTTC	11880
AATTCCTCAT	CACCTGGTTT	TGGATATTTG	CTGGCATCAA	TGACAGGATA	GGCCGCCTGA	11940
TTAACATTGA	TATGOCAGTA	GCCATTGGA	TTTTTCTTGA	GATAGTCTTG	ATGGTAATCC	12000
TCAGCCACCA	CAAAATCTCT	CAAGTTTTC	TTTTCAACTG	CTAGAGGTTG	ATCGTATTTC	12060
TTAGCCAACT	CATCAAGAC	TTGGTTAATC	ACTTCCAAAT	CCTGTGCATC	TGTGTAATAA	12120
ACACCAGTAC	GGTACTGGGT	CCCCACATCA	TTTCCTTGTT	TATTTTGTCT	GGTTGGATTG	12180
ATAATGCGGA	AATAGTGAAG	CAGGATTTC	TTGAGAGAAA	TTTGCTTGCC	ATCATAGGTG	12240
ACATGACGG	TTTCTGCTAG	ACCTGTTTGG	TTAATCAATT	CGTACTTGCT	TGTTTCTCTT	12300
CTACCAATTTG	CATAGCCCTGA	AACGGCATCC	GTACCCCGG	GAACACGTGA	GAAATATTCC	12360
TCCACTCCCC	AGAAACAACC	TCCAGCTAGA	TAAATTTCTG	GCAAGTCTGC	GTCTTTACTA	12420

182

ATTCCTGTTT TTTTCACTGC TTTTCTCTT TGGCTAACTG CCGCCCTTTC AATTGGCAG	12480
GCATCTGTCT GCCCTGCATT TCGTATCAAT AGACATAGA AACCGGTTAT GCGTAGAAAA	12540
AATACTCCTA GCACAAGAA GATTTTAAAC TTATCATCA TAAGACGCCCT CCTAGGCTAA	12600
TTCTCTCAAA GTTTGCAAAA TTGCATCTTT TTCCATGAAT CCGTAGTGTG TTTTGACCAG	12660
CTTGCCCTTCT TTGTCTATA AGGCTTGGGT TGGTAAGAA CGGACACCAT AGTTTCCCAA	12720
AAGTTTCCTT GATGGGTCAA CTAGGACTGG GAGATTTTAA TAATCCAATC CCTTATACCA	12780
ATTCTTAAAG TCGCGTTCAG ATTGCTCTCC CTTATGTCTT GGTGACACTA CTGTCAAGAC	12840
CACATAGTCA TCACCAGCTT CTTTAGCAAT CTCATCCGTA TGTGGAAGAC TAGCCAGACA	12900
GATGGAACAC CAAGAAGCCC AGAATTGTAG ATAGACTTTC TTGCCCTTGT AATCAGATAA	12960
ACGGTAGGTC TTGCATCTA CTCCCATCAA TTCAAAATCA GCCACCTCTT TCCCTTTAGC	13020
TGCGCTGTGT TTAATAGCTG TCTGCTCCGT CTTCAATTCA TCTTTTGGTT GGTGTTCACT	13080
AGTCACGAC TTGCTGAAC AAGCCGTCAA ACAAGGAGC GAACCTGCTC CAGAACACA	13140
TGTTTGCCAT TTTTTCATAT TGATATTCCT TTCCATTTTA TTCAAATAAT TGACTTAAAA	13200
TTGAAGCATT TCCAACAGA ACCAAGAAGC CCATCACAAT AATGAGAJAA CCACCCACTT	13260
TTTTGAGGAT TCCGAGATAG GGATGAAGTT TTCCGAAATG TTTCAAJACA TAACTAGAGG	13320
TCAGAGCTAG AAGCAAGAAT GGTAGCGCCA AGCCCAAGCT ATACACCAAC ATGAGACCAG	13380
CTCCCTGCCA AGCTCTGAA CCACCTGAGC CGGCCAAGGC CAJAJACAGC CCCAGAACC	13440
GCCCCACGCA AGGCGTCCAA GCAAAATCAA AGGTCAAGCC CAATAAAAAA GCCTGACTAT	13500
AGCCCTTACC ATTTTGCCCC TGTCTTGCA GTTGTAGCTT CTTTTCCTTA TAAAGCCCCCT	13560
TAAAGGTAG AATCTCCATT TGGTGCAAC CAAGAAGGAT AATAATTGCC CCAGTAAGAT	13620
ATTGGAACCA AGAAGCATAA AGCAAAATCC CTAAAAAACC AGTCCATAG CCCAACAAAA	13680
TAAATATAAA GGAATTCCT GCTATAAAG CCAGAGTTCG TAATAAATA GTAACGAGA	13740
TTGAAATTT GCGCTAGAA GCCTGAGCAC CATCCTTATC ATCTAGTAAC ACTCCTGTAT	13800
AGACCGTAA CAJAGGTAAG ATACAAGGAG AAJAGAAGGA TAGAATCCCT GCCAAJAGA	13860
CACCTAGAAA AAAGAAAAA TGACCATAA AGTCTCTCT ATCATTTTAT TGATAGATT	13920
ATTATA	13926

(2) INFORMATION FOR SEQ ID No: 6:

- (1) SEQUENCE CHARACTERISTICS:
  - (A) LENGTH: 20199 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 6:

CCCAGCAGAA AANTGGCATT TGGAGATAAT GGAATCGTA AAAAACTAT GTTTGAGAAA	60
ATAACCTTGT TTATCGTGAT TATCATGCTA GTAGCAAGTT TATTGGGAAT TTTTGCAACT	120
GCAATTGGTG CCCTCAGTAA TCTATAAAAT AGATTCAAGA AAATTTAGTG ACTGGGATTT	180
CCCAGCCCTT TTTTAAAGTG AGAAGAAATA ATGAGTATGT TTTTAGATAC AGCTAAGATT	240
AAGGTC AAGG CTGGTAATGG TGGCGATGGT ATGGTTGCCCT TTCGTCTGTA AAAATATGTC	300
CCTAATGGAG GCCCTGGGG TGGTGATGGT GGTCTGGAG GCAATGTGGT CTTCGTTGTA	360
GACGAAGGAC TACGTACCTT GATGGATTTC CGCTACAATC GTCAATTCAA GGCTGAATCT	420
GGTGA AAAAG GGATGACCAA AGGGATGCAT GGTCTGGTG CTGAGGACCT TAGAGTTGGA	480
GTACCA AAG GTACGACTGT TCGTGATGCG GAGACTGGCA AGGTTTTAAC AGATTTGATT	540
GAACATGGGC AAGAATTTAT CGTTGCCAC GGTGGTCTG GTGGACGTGG AAATAATTCGT	600
TTCCGCACAC CAAAAATCC TGCAACCGAA ATCTCTGAAA ATGGAGAACC AGGT CAGGAA	660
CGTGAGTTAC AATTGGAAC TAAAAATCTTG GCAGATGTCG GTTTAGTAGG ATCCCACTCT	720
GTAGGGAAGT CACACCTTTT AAGTGTATT ACCTCAGCTA AGCCTAAAA TGGTGCCATC	780
CACTTTACCA CTATTGTACC AAATTTAGGT ATGGTTGCGA CCCAATCAGG TGAAATCCITT	840
GCAGTAGCCG ACTTGCCAGG TTTGATTGAA GGGGCTAGTC AAGGTGTTGG TTTGGGAAC	900
CAGTTCTCC GTCACATGA GGTACACGT GTTATCCTTC ACATCATTTA TATGTCAGCT	960
AGCGAGGGCC GTGATCCATA TGAGGACTAC CTAGCTATCA ATAAAGAGCT GGAGTCTTAC	1020
AATCTTCGCC TCAATGGAGC TCCACAGATT ATTGTAGCTA ATAAAGATGA CATGCCGTAG	1080
AGTCAGGAAA ATCTTGAAGA CTTTAAGAAA AAATGGCTG AAAATATGA TGAATTTGAA	1140
GAGTTACCA GCTATCTTCCC AATTTCTGGA TTGACCAAGC AAGGTCTGGC AACACTTTTA	1200
GATGTACAG CTGAATGTTT AGACAAGACA CCAGAAATTT TGCTCTACGA CGAGTCCGAT	1260
ATGGAAGAAG AAGCTTACTA TGGATTGAC GAAGAAGAAA AAGCCTTTGA AATTAGTCGT	1320
GATGACGATG CGACATGGGT ACITTTCTGGT GAAAACTCA TGAACCTCTT TAATATGACC	1380
AACCTTGATC GTGATGAATC TGTCAAGAAA TTGCCCCGTC AGCTTCGTGG TATGGGGGTT	1440
GATGAAGCCC TTCGTGCGCG TGGAGCTAAA GATGGGGATT TGGTCCGAT TGGTAAATTT	1500
GAGTTTGAAT TTTGAGACTA GGAGACTGGT ATGGAGATA AACCGATATC TTTCCGAGAT	1560
CGGATGGTA ATTTTGTCTC CGCCGAGAC GTTTGGAATG AAAAGAAATT GGAAGAATA	1620

		184	
TTTAATCGTC	TCAATCCAAA	TCTGCGCTTG	AGATTGGCAC
GAACATAAAA	GGAAAAATCCA	1690	
TCTCAGTAAA	GAAGCTAAAA	AATCCCGTGC	CTCATCAGAC
ACGGGATTTT	GTGCTACGAC	1740	
AGGCATGTAT	AGCAAACTGA	ATCTGGAATA	GCACAGCATA
TCTTCTAAAA	TATAGTAAAA	1800	
TGAAATGAGA	ACAGACAAAA	TCGATCAGGA	CAGTAAAAATC
GATTTCTAAC	AATGTTTTAT	1860	
AACGAGAGAT	GTACTATTCT	AGTTTCAATC	AACTATATTG
TTATAAATTG	ATTTGAATTT	1920	
CAAAATTAAA	TTGTTTGATT	CTTATTTCAA	TTTGTATAG
TATATCTGAT	GTCAAAAGTTC	1980	
TCGGCGAGTC	AAATAGCGAT	TCCCAAGCCT	GAATATCGTG
AGGTAGCGGA	TTAAAAATGGT	2040	
CTGGGGATAG	ACCGTTTTAA	GTCTGACGCT	GGAAATAAGA
ATTTCTCAGAA	GAAGGATAG	2100	
CGAAATCGTG	GCTCTACGAA	CAGGAACGTG	ATAATAAGGC
GTATATAGCG	GATAAGAGGG	2160	
CATCAAACTC	TAAAGTCCAA	AAAGGTAGTC	GTAACTTATA
TGCCTAAATC	ACGAGAGTAA	2220	
TTGAATTCGT	ACTAAGATTT	TCTATTTTCA	CTGTAACTTT
TTAAACGCCCT	TATATCTTGT	2280	
ATACCGAGAG	AAAGATGTAC	GAATATCCCC	GTGAGGTCTA
TCACTATAAA	GAGAAACGA	2340	
CAGATAGAAG	TGATCCTGAG	TCACGGTTAT	CTGTCTGATA
GGACGGTATG	TATAAAACGC	2400	
TTCTGTGAAC	TGAGAGAAGG	GGGAGAAGTT	CTTGCTAAAA
TTTAGTTGAA	CAGCCGTATT	2460	
CCGATACCTA	GATAAGAGAT	CTAGTCTTAG	CTCCTACTCA
GTTTTAGGGG	ATAAAAAAGG	2520	
GGCAATAGCG	ATTTCGAGAA	GATTATACTC	TTTCGAAATC
ACGTCAATAT	2580		
CGCCTTGTCG	TATGTGTAGG	ATACTGACTA	CGTCAGTTCC
ATCTACAACC	TCAAAACAGT	2640	
GTTFPGAGCA	ACCTGCGGCT	AGTTTCCCTAG	TTTGATCTTT
GATTTTCTAT	GAGTATTAGT	2700	
AATTCAGTTA	CTAACTCGTC	AACTCTGATT	TATCCAATAA
AATTGAAAAG	GATGGAAAAA	2760	
AGGATAAATT	TATGATATAC	TTTATTTTGA	AGACCTTATT
AGAAATCTTG	AAAGAGTATT	2820	
GAAAACTTAG	AATGAGAAAA	ATTGTTATCA	ATGTGTGGATT
ACCACTGCAA	GTGAAATCA	2880	
CTATTATGCG	TGCTAAAAAT	AGTGTCTGTT	CCTTAATTCC
AGCTATTATC	TTGGCTGATG	2940	
ATGTGCTGAC	TTTGGATTGC	GTTCAGATA	TTTCGGATGT
AGCCAGTCTT	GTTCGAAATCA	3000	
TGGAATIGAT	GGGAGTACT	GTTAAGCGTT	ATGACGATGT
ATTGAGATTT	GACCCAAGAG	3060	
GTGTTCAAAA	TATTTCAATG	CCTTAATGTA	AAATTAACAG
TCTTCTTGCA	TCTTACTATT	3120	
TTTATGGGAG	CCTCTTAGGC	CGTTTGGGTG	AAGCGACAGT
TGGTCTACCG	GGAGGATGTG	3180	
ATCTTGTGTC	TGTTCCGATT	GAATTAACCC	TTAAGCGGTT
TGAAGCTATG	GGTGCCACTG	3240	
CTAGTACAGA	GGGAGATAAC	ATGAAGTTAT	CTGCTAAAGA
TACAGGACTT	CATGGTCCAA	3300	
GTATTTACAT	GGATACGCTT	AGTGTGGGAG	CAACGATTAA
TACGATGATT	GCTGCCGTTA	3360	
AAOCAAATGG	TGCTACTATT	ATTGAANAATG	CAGCCCGTGA
ACCTGAGATT	ATTGATGTAG	3420	

CTACTCTCTT GAATAATATG GGTGCCATA TCCGTGGGCG AGCAACTAAT ATCATCATTA 3480  
 TTGATGGTGT TGAAGATTA CATGGGACAC GTTCATCAGGT GATTCCAGAC CGCATTTGAAG 3540  
 CTGGAACATA TATATCTTTA GCTGCTGCAG TTGCTAAAGG AATTCGTATA AATAATGTTT 3600  
 TTTACGAACA CCTGGAGGG TTTATTTGCTA AGTTGGAAGA AATGGGAGTG AGAATGACTG 3660  
 TATCTGAAGA CAGCATTTTT GTGAGGGAAC AGTCTAATTT GAAAGCAATC AATATTAAAG 3720  
 CAGCTCTCTA CCCAGGCTTT GCAACTGATT TGCAACAACC GCTTACCCCT CTTTACTTAA 3780  
 GAGCGAATGG TCGTGGTACA ATTGTCGATA CGATTACGA AAAACGTGTA AATCATGTTT 3840  
 TTGAACTAGC AAGATGGAT GCGGATATTT CGACAACAA TGGTCATATT TTGTACACGG 3900  
 GTGGACGTGA TTTACCTGGG GCCAGTGTTA AAGCGACCGA CTTAAGAGCT GGGGCTGCAC 3960  
 TAGTCATTGC TGGGCTTATG GCTGAAGTA AAACCTGAAAT TACCAATATC GAGTTTATCT 4020  
 TACCTGGTTA TTCTGATATT ATCGAAAAAT TACGTAATTT AGGAGCGAT ATTAGACTTG 4080  
 TTGAGGATTA AACCGTAGAG GTGTTTATGA ATATTGGAC CAAATAGCA ATGTTTCTCT 4140  
 TTTTGAAC GGATCGCTTG TATTTCGCTT CTTTCTTTT TAGTGATAGT CAGGACTTCC 4200  
 GCGAGATAGC TTCAAAATCCA GAAATCTTC AATTATTTT CCCAACGCGAG GCAAGCTGCG 4260  
 AAGAAAGTCA ATATGCACTG GCCAATTACT TTATGAAGT CCCTTTGGGA GTGTGGGCAA 4320  
 TTTGTGACCA GAAAAATCAA CAAATGATTG GTTCTATTAA ATTTGAGAAG TTAGATGAAA 4380  
 TCAAAAAAGA AGCTGAGCTT GGCTATTTTT TGAGAAAAG TGCTTAGTGC CAAGGATTTA 4440  
 TGACAGAGGT TGTTAGAAAA ATTTGTCAGC TTTCTTTTGA GGAATTTGGC TTAAAAAAT 4500  
 TATTATTATCAT TACCCACCTT GAAAAATAAG CTAGCCAAAG AGTTGCTCTT AAGTCTGGAT 4560  
 TTAGTTTGTG CCGTCAGTTT AAGGGAAGTG ATCGTTACAC AAGAAAAATG CGGGATTATC 4620  
 TTGAATTTTG GTATGTAAAA GGAGAGTTCA ATGAGTAAGC ATCAGGAAAT TCTAAGCTAT 4680  
 TTGGAGGAAT TACCAGTAGG TAAAGGGTC AGTGTTCGTA GCATTTGGA TCACTAGGA 4740  
 GTTAGTGATG GAACAGCCTA TCGGGCTATT AAAGAAGCTG AAAACCTGTG AATTGTGGAG 4800  
 ACCGCTCTTA GAAGTGAAC AATTCGTGTT AAATCCAGA AAGTTGCTAT AGAGAGATTA 4860  
 ACGTTTGCTG AAATTCGAGA AGTGACTTCT TCTGAGGTTT TGGCTGGCA AGAAGGTTTA 4920  
 GAGAGAGAA TTAGTAAGTT TTCAATTGGT GCCATGACTG AACAAAATAT CTTGTCTTAC 4980  
 CTTCTGATG GGGGGCTCTT GATTGTGGA GACCGAACCC GTATTAGTT GCTAGCCTTG 5040  
 GAAAATGAAA ATGCACTTCT GGTACAGGG GGATTTCAGG TTCATGATGA TGTGCTTAAA 5100  
 CTGGCCAAAT AAAAAGGAT TCTGTCTTA AGAAGTAAGC ATGATACCTT TACCGTCGG 5160

186

ACCATGATCA ATAAAGCCCT	GTCAAATGTC CAAATCAAGA	CTGATATTCT GACAGTTGAG	5220
AAACTTTATC GCCCTAGTCA	TGAGTATGOT TTCTCGAGAG	AGACAGATAC AGTTAAAGAT	5230
TATTGGGACT TGGTCTGTAA	GAATCGTAGC AGCCGTTTCC	CTGTTATCAA TCAACATCAG	5340
GTGCTGTGTG GTGTTGTAA	CATGAGAGAC GCTGGTGATA	AATCACCAGG CACGACAATT	5400
GATAGGTTA TGTCTGTAG	TCTATTTTGG GTTGSATTAT	CGACAAATAT TGCCAATGTG	5460
AGTCAACGGA TGATCGCAGA	AGACTTTGAA ATGGTACCAG	TTGTTGGAAG CAATCAAATC	5520
TTGCTTGGCG TTGTGACGCG	ACGAGATGTC ATGGAGAAGA	TGAGCCGCTC CCAAGTTTGG	5580
GCTCTACCAA CTTTTCTGTA	GCAGATTGGA CAAAAGCTCT	CTTATCACCA TGATGAAGTA	5640
GTCAATTACAG TGGAACTCTT	TATGCTAGAA AAAAATGGAG	TTTTGGCTAA TGGTGTATTG	5700
GCAGAAATTC TGACCCACAT	GACCCGATTT AGTTGTTAAT	AGTGTCTGCA ATCTCATTAT	5760
CGAGCAGATG CTGACTACT	TTTTGACGCG TGTTACAGATA	GATGATATAT TGCGCATTCA	5820
GGCAGCGATT ATTCATCATA	CGACACGGTC AGCTATAAAT	GATTACGATA TTTATCATGG	5880
TCACCAAGATT GTTTCAAAG	CAAAATGTGAC TGTAAAAATT	AATTAGAAAC TAGGAGAAAA	5940
GATGATAACA TTAJAATCAG	CTCGTGAATC CGAAGCTATG	GACAAGGCTG GTGATTTTCT	6000
AGCAAGTATT CATATAGGCT	TACGTGATTT GATTAAGCCA	GGCCTAGATA TGTGGGAAGT	6060
TGAGAATAT GTCCGCGCTC	GTGTGTAAGA AGAAAAATTC	CTTCCACTTC AGATTGGGGT	6120
TGACCGTGCC ATGATGGACT	ATCCTTATGC TACCTGTGTC	TCTCTTAAAG ATGAAGTGGC	6180
TCACCGTTTC CCTCGTCATT	ATATCTTGAA AGATGGTGAT	TTGCTCAAAG TTGATATGGT	6240
TTTGGGAGGT CCCATTGCTA	AATCTGACCT AAAATGCTCA	AAATTAACT TCAACAATGT	6300
TGAACAAATG AAAAAATACA	CTCAGAGCTA TTCTGTGTGT	TTAGCAGACT CATGTGGGGC	6360
TTATGCTGTT GGTACACCGT	CGAAGAAGT CAAAACTTG	ATGGATGTAA CCAAGAAGC	6420
TATGTACAAG GGTATTGAGC	AAGCTGTGTG TGGAAATCGT	ATCGGTGATA TCGGTGCGGC	6480
TATTCAAGAA TACGCTGAAA	GTGCTGGTTA CGGTGTATGT	CGTGATTGGG TTGGTCATGG	6540
TGTTGGCCCA ACTATGCACG	AAGAACCAAT GGTTCCTAAC	TATGTATTATG CAGGTCGTGG	6600
ACTCCGCTCT CGTGAAGGAA	TGGTCTTAAC CATTGAACCA	ATGATCAATA CAGGCGATTG	6660
GGAAATTGAT ACAGATATGA	AACTGGTTG GCGCATTAAG	ACCATTGACG GTGATTGTCT	6720
ATGTCAGTAT GAACACCAAT	TTGCTATTAC GAAAGATGGA	CCTGTATATC TGACTAGCCA	6780
AGGTGAAGAA GGAACCTATT	AATAAAAAGT GAAAGACTA	CTGGAAGTTT ATTTTGATAA	6840
AAAAATCAAT AGATCTTTTC	ATPAATAAAC GCATTGTATC	AAAGTGTAGG GGTGATATC	6900
ATCGCTTTT CTGCTTTTAA	GATTTTTCCT AACTCTGTTT	GTAAGCGCAT CATAACAAAG	6960



GGTCTAGGAT TCAGGGCTCT	CCTCCTATAT ACTATTAGTA AAGTAAACT AAGGGAGGAT	7020
ATTTTAGTGT CGCAGTCTAT	TGTTCCTGTA GAGATTCCAC AATATTGTCG TTTTGATTCT	7080
AAAAAGAGAA ATGGAATCT	GTTTAATGTT CGTATTGCCA ATCTTAAAT TACTTTTTTA	7140
TATTATACTT CTTCGCAAC	AAAATATGGT ATAGTAGTTC TATGAATGAT GAGCAAGTA	7200
AACAACATAA TGATGCACGA	TTTAAGCGTC TTGTTGGTGT TCAGCGTACC ACTTTTGAAG	7260
AGATGTTAGC TGTATTAAAA	ACAGCTTATC AACTTAAACA CGCAAAAGGT GACGAAAC	7320
CTAAATTAGC CCTAGAAGAC	CTTCTTATGC CCACTCTTCA ATAGTGCAG AATATCGAAC	7380
TTATGAAGAA ATTGCGGCTG	ATTTTGGTAT TCACGAAAGC AACTTTATCC GTCGGAOCCA	7440
ATGGGTGAA ATAACTCTTG	TTCAAAGTGG TTTTACGGTT TCAAGAACTC CTCTCAGTTC	7500
TGAGGACACG GTAATGATTG	ATGCGACGGA AGTAAAAATC AATCGCCCTA AAAAACAAT	7560
TAGCGAATGA TTCTGGTAAA	AAGAAATTTT ACGCTATGAA GGCTCAAGCG ATTGTCAACA	7620
GTCAAGGGAG AATTGTTTCT	TTGGATATCG CTGTGAAC TAAGTATGAT ATGAAGTTGT	7680
TCAAAATGAG TCGTAGAAAT	ATCGAACAG CTGGTAAATC CTGGCTGAC AGTGGTTATC	7740
AAGGGCTCAT GAAGATATAT	CCTCAAGCAC AAACCTCCAG TAAATCCAGC AAACCTCAAGC	7800
CGCTAACAGC TGAAGATAAA	GCCTATAACC ATGCGCTATC TAAGGAAGA AGCAAGGTTG	7860
AGAACATCTT TGCCAAAGTA	AAAACGTTTA AATATTTTTC AACACCTATC CGAAATCATC	7920
GTAAACGCTT CGGATTACGA	ATGAATTTGA GTGCTGGTAT TATCAATCAT GAACAGGAT	7980
TCTAGTTTGG CAGGAAGTCT	ATTGAGGTAT TGAGCTAGTT TATGAAAAA TTGGGTGAAA	8040
AGTCAGTGT TTTAGAAACC	CACAGTGTAG TATTTAGTT TCAATCCACT ATATTTTGTCT	8100
ACTCCCCGTA AAGTTTCTAT	TTTCCCTGAT TTCTGATATA ATAGAAATAT TGACTTCAAG	8160
AGTAAGGAAG AGAAGATGAA	CGCATTATTA AATGGAATGA ATGACCGTCA GGCTGAGCGC	8220
GTGCAACGA CAGAAGGTCC	CTTGCTAATC ATGGCAGGGG CTGGTCTGG AAAGACTCGT	8280
GTTTGAAGCC ACGTAGTCC	TTATTGATT GATGAAAAGC TGCTCAATCC TTGGAATATC	8340
TTGCCATTA CCTTTACCAA	CAAGGCTGCG CGTGAGATGA AAGAGCGTGC TTATAGCCTC	8400
AATCCAGCGA CTCAGGACTG	TCTGATTGCG ACCTTCCACT CCAATGTGTG GCGTATTTTG	8460
CGTCGCGATG CGGACCATAT	TGGCTACAAT CGTAATTTTA CAATGTGGA TCCTGGTGAA	8520
CAGCGAAGCG TCATGAACG	TATTTCTCAA CAGTTGAATC TGGACCTTAA AAAATGGAAT	8580
GAACGAACATA TTTTGGGGAC	CAITTTCCAAT GCTAAGAATG ATTTGATTGA TGATGTTGCT	8640
TATGCTGCCC AAGCTGGCGA	TATGTATACG CAAATGTGG CCGAGTGTTA TACAGCTAT	8700

188

CAAAAGAAC TTGTCAGTC TGAATCCGT GACTTTGATG ATTTGATTAT GCTGACCTTG	8760
CGTCTCTTTG ATCAAAATCC TGATGTTTGG ACCTACTACC AGCAAAATTT CCAATACATC	8820
CACGTTGATG AGTACCAAGA TACCAACCAC GCTCAGTACC AATTGGTCAA ACTCTTGGCT	8880
TCCTGCTTAA AAAATATCTG TGTGGTTGGG GATGCGGACC AGTCTATCTA CGGTGGCCT	8940
GGTGTGATA TGCAGATAT CTGGACTTT GAAAGGATT ACCCCAAGC CAGGTGTGTT	9000
TTGTTGGAGG AAAATTACCG CTCACCAAAA ACCATTCTCC AAGCGGCCAA CGAGGTTAT	9060
AAAAATATA AAAATCCCG TCCTAAAAAT CTCTGGACTC AAAACGCTGA TGGGAGCAA	9120
ATCGTTTACT ATCGTGCGA TGATGAGCTG GATGAGGCTG TATTTGTAGC CAGAACCATC	9180
GATGAACCTA TCGCAGTCA AAACCTCCTT CATTAAGATT TTGCAGTTCT CTATCGGACT	9240
AATGCCAGT CCGGTACAAT TGAGGAAGCC CTGCTCAAGT CTACATTTCC TTATACCATG	9300
GTGGCGGAA CCAATTTCTA CAGCCGTAAG GAAATTCGCG ATATTATTGC TTATCTCAAC	9360
CTTATTGCTA ATTTGAATGA CAATATTAGT TTTGAGCTA TTATCAACGA GCCTAAACGT	9420
GGAATTGGTC TAGGTACAGT TGAGAAAATC CGTGATTTTG CAAATTTGCA AAATATGTCT	9480
ATGCTGGATG CTTCTGCTAA TATTATGTTG TCTGGTATCA AGGGTAAGGC AGCCCAATCT	9540
ATCTGGGATT TTGCAATAT GATGCTTGAT TTGCGGGAGC AGCTAGACCA CTTAAGCAAT	9600
ACAGAGTTGG TTGAGTCCGT CCTAGAAAA ACAGGTATAG TCGATATTCT TAACCTCCAA	9660
GCGACTCTAG AAGCAAGGC ACGGGTTGAA AATATCGAAG AGTTTCTTTC TGTACGAAG	9720
AACTTTGATG ACACCACGGA TGTGACAGAA GAGGAAACTG GTCTGGACAA ACTGAGTCGT	9780
TTCTTAAATG ACTTGGCTTT GATTGCCGAC ACAGATTACG TAGTCAGGA GACATCAGAA	9840
GTGACCTTGA TGACCCCTGA TGCTGCCAAA GGTCTCGAAT TTCCAGTTGT CTTTTGTATT	9900
GGGATGGAAG AAAATGTCTT TCCACTTAGT COTGCGACTG AAGATTGAGA TGAATTAGAA	9960
GAAGAGCGCC GTCTAGCCTA TGTAGGTATC ACGCGTGACG AGAAAATCTT CTATCTGACC	10020
AATGCCAACT CACGCTTGCT TTTTGGTCTG ACCAATATA ACCGTCCGAC TCGTPTTATT	10080
AACGAATCA GTTCAGACTT GCTTGAATAT CAAGGTCTGG CTCGTCTGCA AAATACAGC	10140
TTTAAGGCAT CATATAGCAG TGGTAGTATT TCCTTTGCTG AAGGTATGAG TTTGCTTCAG	10200
CTCTTTCAAG ACCGTAAACG CGGTGCTGCC CCAAAATCAA TCCAGTCAAG CGGTCTTCCA	10260
TTTGGTCAAT TTACAGCTGG CGCAAAACCA GCATCTAGCG AGGCAAAATT GTCCATTGAT	10320
GATATTGCTC TCCACAGAA ATGGGGAGAG GGAACCGTTC TGGAAATTTC AGGTAGCGGT	10380
GCTAGGCAGG AATTGAAAA CAATTTCCCA GAAGTAGGTT TGAAAAAATC TTTAGCCAGT	10440
GTGCTCCAA TTGAGAAAAA AATCTAATTT TCCATCTTTC TCACGAATAA TAAAGTGAGG	10500

AGGATTTTTA TGTACAGTAT TTCATTCCAA GAAGATTAC TATTACCAAG AGAAAGGCTG 10560  
 OCCAAGGAAG GAGTTGAAGC GCTTAGTAAC CAAGAGTTGC TAGCTATTTT ACTCAGGACA 10620  
 GGAACACGTC AAGCTAGCGT TTTTGAAATT GCCCAAAAG TCTTGAACAA TCTTTCAAGC 10680  
 CTAACGAGTT TGAAAAAAT GACCCGCGAG GAAATTGCAGA GTTCTCTGG TATTGGGCGT 10740  
 GTTAAGGCCA TAGAATTACA AGCTATGATT GAATGGGGC ATCGTATTCA CAAACAGGAG 10800  
 ACTCTTGAAA TGGAAAGTAT TCTCAGCAAT CAAAGTTGG CCAAGAGAT GCAGCAGGAA 10860  
 TTAGGGGATA AAAACAAGA GCACCTGGTG GCACCTCTATC TCAATCTCA AATCAAATC 10920  
 ATCCATCAGC AGACCATTTT TATCGGGTCT GTAACCTGTA GTATCGCTGA ACCGCGAGAG 10980  
 ATTCTTCACT ATGCAMTCAA GCATATGGGC ACTTCTCTTA TCTTGCTCCA CAATCATCCT 11040  
 TCAGGAGCGG TAGCGCTAG CCAAAATGAT GATCATGTCA CTAAACTTGT TAAAGAAGCC 11100  
 TCGCAATTGA TGGGATTGTT TCTCTTGGAC CATTTGATTG TCTCTCATTC TAATTACTTT 11160  
 AGTTATCGTG AAAAGACAGA TTTAATCTAA AGTTCATTAA CGACATAGTC AAAGAGTTT 11220  
 TTATCTTTGG GACGATTTTC AAAAAGAAAT TCTGGATGCC ATTGGACACC GAGAAAGGG 11280  
 ACATCATCCG TACTCATGAC AGCCTCAAT ATACCATCTT TAGGATCATG AGCCACAAC 11340  
 TTTAAATTG GTGCTAAGTC CTGTATGCTC TGGTGGTGA AGGAGTTGAT ATGAGAGATT 11400  
 TCTCCATAGA TTTCTTGGG AACGGTATCT GGTCTGTATA CCAAGCGTTG AGTTGTGTAC 11460  
 TCAACAGAAG AATCCTGCCA ATGGTCTTGG ATATCTTGGT ACAAAGTCC ACCCATGGCA 11520  
 ACGTTAAAGA GTTGGGTACC ACGGCAGACA GAGAAAAATG GCTTTTTCTG TTTAATAGCT 11580  
 TCCTTGATGA GGGCCAGTTT GAAGATATCT CTTTGAAGGT GATAGTCATC ACTATCAATG 11640  
 GTTTTGGGTT CGCCATAAAA TTTTGGATCG ACATTTTGGC CACCTGTCAA GATGAGCTTG 11700  
 TCAATCAAACT TGATATAGTG GCAGGCCATT TCTTGATCAC CAATCGTAG GATGATGGGA 11760  
 ATCCCTCCAG CATCTTTAAC GCCTTCAACA AAGCCTTTTG CTGCGTAGCT CATCATGATG 11820  
 TCAATCATCTG GATGAGTTT TTCGTTTCTT GTAATCCCAA TAACTGGTTT TTTCATAAAA 11880  
 TGATTTTCGC TTTCTAATCC TCTTTTCGCA TGAAGTAGAG GAGGGTTGG AGTTCACTTG 11940  
 TCAATCGAC ATACTGACG ACCACGTCCT TGGTAAATG CAGATGGACT GGTGAAAAAC 12000  
 TGAGAAATCC TTTACACCA GCATCAACCA AGAGATTAGC AACCTCTGT GACTTGACGC 12060  
 TGGGAACAGT TAGGATAGCA GTCTTCACAT CAGCATCCTT GATTTTATCC TTGATCTGAG 12120  
 AAATCCCGTA AATGGGAATC CCGTCAGGAG TTTGGGTACC GACTTCAAGA TGGTCGTCTA 12180  
 GGTCAAAGGC CATGATAATC TTAATCTTGT TACGTTCTGT GAAGCGGTAG TGGAGAAGGG 12240

190

CATGGCCCA	ATTTC	CAAC	CGCA	TGACAT	GTG	TCATT	12300
AATCGG	AAAT	GT	AGTT	TTT	GT	CAAT	12360
CACCA	AAAT	CA	CGAC	TAC	TC	CT	12420
TTT	GT	CT	TAG	AT	TC	AT	12480
AGAGAG	AG	AG	AG	AG	AG	AG	12540
CACA	AC	CT	TT	CT	TT	CT	12600
AAAA	CT	AA	AG	CT	CT	AG	12660
AG	GT	CT	CC	GA	AG	CT	12720
AG	GT	AT	CT	CT	GA	AG	12780
TAC	TT	GA	AG	AG	AG	AG	12840
CT	AT	CT	TT	CT	GA	AG	12900
TTTT	TT	CT	GA	AG	AG	AG	12960
CT	GT	CA	GT	AG	AT	TT	13020
GA	AG	CT	TT	CA	AG	AG	13080
CT	TA	GT	TT	CA	AG	AG	13140
GT	AA	TC	GA	AT	TT	CA	13200
CG	AT	AT	TT	TT	GA	AG	13260
AC	AA	CT	GA	AG	AG	AG	13320
GA	AC	AT	GG	CT	GA	AG	13380
CAG	CA	GG	AG	AG	AG	AG	13440
GA	AC	ATA	AG	AG	AG	AG	13500
CC	AG	AA	AG	AG	AG	AG	13560
CT	CT	TT	AG	AG	AG	AG	13620
TT	TT	CT	GT	AG	AG	AG	13680
GA	AA	CC	GA	AG	AG	AG	13740
CCA	AT	CT	TC	GA	AG	AG	13800
CAG	CA	GG	AG	AG	AG	AG	13860
CAT	TAT	TG	AC	AG	AG	AG	13920
GCT	CG	CG	CT	GT	AG	AG	13980
GA	CC	CA	AG	AG	AG	AG	14040

GTTCCTCTTG	TTCTTGGTGA	CGAAGACAGT	AGCCAAATGAT	GGTAGTATTA	TTGCCCTCAG	14100
TCCACACAGA	AGTGAAAAAG	ATATGTGTGAG	GTTCCTGTCT	TAGTAACCTGG	GCTAGTTCCCT	14160
GACGGGCTTC	TCGCAGAGGT	TTGCCAGCTT	GAGGACCATG	ACCATGAATA	CTAGAAGGAT	14220
TTCCGTGGGT	TTCTTGCAATA	ACCTTGGTCA	TAGCTGAAAT	AGCAACTGCT	GACATAGGAG	14280
TCCTTTCGAGC	ATTGTCCAAA	TAAATCAAG	AATCACCTTA	TTCTCTTTTA	TTGTAGGCAG	14340
AGATGTGGCT	GACTGGTTTT	CTTTCCGTGA	TACGGACGAT	AGCATCACCA	ATTAACTCAC	14400
TAGCAGTGAT	GTAGCATACA	TTTTTAGGAG	TTTTTTCTTT	TGTTGCTACT	GAATCAGTCA	14460
CAGAAATTC	TTTAATATTA	GTATTGTCAA	GAGCTCAGC	AGCTCCCTCG	ACGAAGAGAC	14520
CGTGGCTAGA	AACAGCATAA	ATTTCTGTAG	CTCCTTCAG	TTCAACGATT	TTAGAAGCTT	14580
CAGAGAAGGT	ACGTCTGTGA	TTTAAAAATAT	CATCAATCAA	GATAGCTTTC	TTACCTTCAG	14640
CATCACCAAT	AAATAAACCT	TCGTTACGAG	TTGCATCTTC	TTGAGGGTAG	TCGATAATGG	14700
CGATAGGAGC	ATCAAGATAT	TCAGCCAGGC	TACGCGCAGC	TTTGACACCT	GAMPTTTTAG	14760
GGCTAACGAC	AACAACATCT	GAACCAAGCA	ATCCTTTATC	GCAATTAATG	TTTGCGAATA	14820
GGGGAACAGT	GAAAAGATTA	TCCACTGGAA	TATCAAGAA	ACCTTGAACC	TGAACGGCAT	14880
GCAAAATCAAG	AGTCAGGATA	CGATCAACTC	CAGCTTTAAC	CAGCATAATTG	GCAACTAGTT	14940
TTGTCTTTAAG	TGGCTCACGA	GGACAAAGCA	TGCGGTCTTG	ACGTGCATAG	CCAAAATATG	15000
GAAGSACAAC	GTGTACTACTG	TGGGCACTTG	CACGCACACA	AGCATCGACC	ATGATTTAACA	15060
ATTCACATTAG	GTGGTTGTTG	ACAGGGAAAC	TTGTGTGATTG	GATGATGTAA	ACATCATAAC	15120
CACGACACT	TTCTTCGATA	TTTACTTTGGA	TTCTCTCCGC	TGAAAATTGA	CGTATGATA	15180
GTTTTCACAG	TGGGACACCA	ACAGCTTGGG	CAATTTTGTG	TGCAATCTCT	TGGTTAGAGT	15240
TGAGTGCAG	AAGTTTCTATG	TTTTTTCTAT	CTGACATTAT	AGACCGTCT	CTGTAAACTT	15300
TATATACTCT	AGTTATATTT	ACCTTACATA	TATGAATGG	GATTGTGTGA	TTTTTATCTT	15360
TTCTATTTTAA	CCAAAAAATG	GAGATATATT	CAGCTATTTT	TCATACTTTT	GACAAATCGA	15420
ACCAATTTTG	AAGGAGCTTT	TTGATAGGAA	ATCTGATTTT	TCTCTAAAAA	TTCTCGAAAA	15480
TCCTGTTTTG	CTTGCTCATG	ATTTTCCACT	TCAAGCTCCA	ATTGCTAATC	TGTTATATCA	15540
AAGTATCGGC	TCGTATCCAG	TGCCATGAGA	CCAATAGCTG	TTTTCATTTT	ATAGCGAAGC	15600
GTGTGTAGAC	AACCAAGAAC	CTGCCAGTTC	TTACTTTTGA	TACCATTGTT	CGCCAAATTCA	15660
TCCAGTACTA	GCCCTTTAGG	AAGTCTTCC	TTACTCAGAT	AGTTCTCAGC	ATCTTTTAGT	15720
TGCAATTTTT	GGTTGTATTC	CATGTTTCCA	ACACTCTGCG	GGACTTTGAG	TGTCAACTCA	15780

192

GCUCAGTCTT CAAGGTTTCG AATGCGCATA GCGACTTTCT TTCTCGCAG TTCAAAATCA	15840
GGCGTGTGCA TGTAGTAATT TGTTTGAAGA ACAGGAGTGA CACCTGTGAA CTGCTCTTTT	15900
AGACGATTGT ATTCATCTTT TTTCATAGT GTTTTCAATT CAATTTCTAA ATGTTTCATT	15960
TTTCTTACCT TTTTTCATCG TTGAAAGCGG ATTTATGGTA TAATAAGCAT TGTATTTATT	16020
GTATATGAAT CTGGAGAAAA AATCAAGAT ATTTTTCAGC GATAATATGA GAACAAGGGA	16080
GAATATATGA OCTTAGAATG GGAAGAATTT CTAGATCCTT ACATTCAGC TGTGTGTGAG	16140
TTAAAGATTA AACTTCGTGG TATTCGTAAG CAATATCGTA AGCAAAATTA GCATTCCTCCA	16200
ATTGAGTTTG TGACCGGTTCG AGTCAAGCCA ATTGAGAGCA TCAAAAGAAA AATGGCTCGT	16260
CGTGGCATTG CTTATGCGAC CTGGAACAC GATTTCAGG ATATTGCTGG CTTACGTGTG	16320
ATCGTTCAGT TTGTAGATGA CGTCAAGGAA GTAGTGGATA TTTTGCACAA GCGTCAGGAT	16380
ATCGCAATCA TACAGGAGCG AGATTACATT ACTCATAGAA AAGCATAAGC CTATCGTCC	16440
TATCATGTGG TAGTAGAATA TACGGTTGAT ACCATCAATG GAGCTAAGC TATTTTGGCA	16500
GAATTCACAA TTCGTACTTT GSCCATGAAT TTCTGGGCAA CGATAGAACA TTCTCTCAAC	16560
TACAAATACC AAGGGGATT CCAGATGAG ATTAAGAAGC GACTGGAAAT TACAGCTAGA	16620
ATCGCCCATC AGTTGGATGA AGAAATGGGT GAAATTCGTG ATGATATCCA AGAAGCCCAG	16680
GCACCTTTTG ATCCTTTTGA TAGAAAATTA AATGACGGTG TAGGAAACAG TGACGATACA	16740
GATGAAGAAAT ACAGGTAAC GAATTGATCT GATAGCCAAT AGAAAACCGC AGAGTCAAAG	16800
GGTTTGTGAT GAATTGCGAG ATCGTTTGAA GAGAAATCAG TTTATACTCA ATGATACCAA	16860
TCCGATATT GTCAATTCOA TTGGCGGGGA TGGTATGCTC TTGTCGGCCT TTCATAAGTA	16920
CGAAAATCAG CTTGACAAGG TCCGCTTTAT CGGTCTTCAT ACTGGACATT TGGGCTTCTA	16980
TACAGATTAT CGTGATTTTG AGTTGGACAA GCTAGTGACT AATTTCGAGC TAGATACTGG	17040
GGCAAGGGTT TCTTACCCTG TTCTGAATGT GAAGGTCTTT CTTGAAAATG GTGAAGTTAA	17100
GATTTTCAGA GCACTCAACG AAGCCAGCAT CCGCAGGTCT GATCGAACCA TGGTGGCAGA	17160
TATGTATATA AATGGTGTTC CTTTGAACG TTTTCGTGGA GACGGGTAA CAGTTTCGAC	17220
ACCGACTGGT AGTACTGCCT ATAACAAGTC TCTTGGCGGT GGTGTTTAC ACCCTACCAT	17280
TGAAGCTTTG CAATTAAAGG AAATTCGCCG CCTTAATAAT CGTGTCTATC GAACACTGGG	17340
CTCTTCCATT ATTGTGCGTA AGAAGGATAA GATTGAACCT ATTCCACAAA GAACGATTA	17400
TCATFACTATT TCGGTTGACA ATAGCGTTTA TTCTTCCGT AATATTGAGC GTATTGAGTA	17460
TCAAAATCGAC CATCATAGA TTCACTTTGT CGCGACTCCT AGCCATACCA GTTTCGGAA	17520
CCGTGTTAAG GACGCCTTTA TCGCGAGGT GGATGAATGA GTTTGAATT TATCGCAGAT	17580

GAACATCTCA AGGTTAAGAC CTCTTAAAA AAGCAGGAGG TTCTTAAGG ATTGCTGGCC 17640  
 AAGATTAAAT TTCGAGGTGG AGCTATTCTG GTCAATAATC AACCGCAAAA TGCACGTAT 17700  
 CTATTGGACG TTGGAGACTA CGTTACCAIT GACATTCCCG CTGAGAAAGG CTTTGAACCC 17760  
 TTGGAGGCTA TTGAAGCTCC ATTAGATATT CTCTATGAGG ATGACCACCT TCTAGCTTG 17820  
 AATAAACCTT ATGGAGTGGC TTCTATTCTT AGTGTCATC ACTCTAATAC CATTGCCAAT 17880  
 TTATCAAGG GTTACTATGT CAAGCAAAAT TATGAAAATC AGCAGGTTCA CATTGTTACC 17940  
 AGACTAGATA GGGATACTTC TGGCTTGATG CTCTTTGCCA AGCAGGTTA TGCCCATGCA 18000  
 CGATTAGACA AGCAGTTGCA GAAGAAATCT ATCGAGAAAC GCTACTTTGC TTTGGTTAAG 18060  
 GGAGATGACG ATTTGGAGCC AGAAGGGGAA ATTATTGCTC CGATTGGCGG TGATGAAGAT 18120  
 TCCATTATTA CCAGACGAGT GGCTAAAGGC GGAAAGTATG CCCATACTTC ATACAAGATT 18180  
 GTAGCTTCTT ATGGAATAT TCACTTGGTC TATATTCAAC TGCACCTGG TCGAACCCAT 18240  
 CAATCCGAG TCCATTTTTC TCATATCGGT TTCTCTTGC TGGGAGATGA TTTGTATGGT 18300  
 GGTAGTCTGG AAGATGCTAT TCAAGCTCAG GCTCTGCAIT GCCATTACCT ATCTTTTAT 18360  
 CATCCATTTT TAGAGCAAGA CTTGCAAGTA GAAAGTCCCT TGCCGATGA TTTTAGTAAC 18420  
 CTIATTACCC AGTTATCAAC TAATACTCTA TAAAACTGT CTCAGAGTAT AATTATTATC 18480  
 TTAAAGGAGA AAACCTATGG AAGTTTTTGA AGTCTCAAA GCCAACCTTG TTGTTAAAAA 18540  
 TGCTCGTATC GTTCTCCCTG AAGGGGAAGA GCTCGTATT CTCAAGCAA CAAAACGCTT 18600  
 AGTAAAGAA ACAGAAGTGA TTCTGTTTTT GCTTGAATC CTGAAAAA TTAATAATTA 18660  
 TCTTGAATTT GAAGGAATCA TGGATGGTTA TGAGGTATC GACCTCAAC ATTATCCTCA 18720  
 ATTTGAAGAA ATGGTTTCTG CCTTGGTGA GCGTCGAAG GGCAAAATGA CTGAAGAAGA 18780  
 TGTACGCAAG GTTTTGGTTG AAGATGTCAA CTACTTTGGT GTGATGTGG TTTACTTGGG 18840  
 CTTGGTTGAT GGAATGGTGT CAGGAGCGAT TCACTCAACA GCTTCAACAG TTCCGCCAGC 18900  
 TCTCAAAATC ATCAAACTC TCCCAATGT AACTCTGACT TCAGGAGCTT TCCTCATGGT 18960  
 TCGTGTACG GAACGTTACC TATTGGAGA CTGTGCCATT AACATCAATC CAGATGCAGA 19020  
 AGCCTTGGCT GAAATTGCCA TCACTCAGC AATCAAGCT AGATGTTTG GCATCGAACC 19080  
 TAAAAATTGC ATGTTGAGCT ATTCTACTAA AGGTTCAAGG TTTGTTGAAA GCGTTGATTA 19140  
 GGTGTTGAA GCACTAAAA TTGCTCAGCA CTTGGCTCCT GACCTTGAAA TCGATGTTGA 19200  
 GTTGCAATTT GATGCAAGCT TTGTTCTGTA AACTCAGCT CTGAAGCTC CTGGAAGTAC 19260  
 GGTAGCTGGT CAAGCAATG TCTTCATCTT CCCAGTATC GAGGAGGAA ATATTGGTTA 19320

194

CAMGATGGCT GAACGCCTGG GTGGCTTTGC GGCTGTAGGA CCTGTTTTGC AAGGTTTTAA	19380
CAAGCCAGTT AATGATCTTT CTCGTGGATG TAATGCAGAT GATGTTTACA AGTTGACCTT	19440
CATCACAGCA GCTCAAGCAG TTCAATCAATA GTGAAACTA TAAAGTGATA TACTATGCTA	19500
TACTGTAGTT ATGAACCTAT GTACGAAAG CACTGCCATT AATTCCTGAG AACTAAATTA	19560
CTGATTTGGT TCAAAAAGGA AAACCTCCAA GCGATGATAT CCTGTCTATA CACGACCTAT	19620
AGAAATCTGT AATATACATA TCCGTAAAC GATAAATTC CTTTTTGATT TTAATGAGT	19680
ATGAAAAGAG AATTTTTTGG CTCTTTGTCA ACTGTAGTG GTTGAAGAAA AGCTAAGCTC	19740
GAGAAAGGAC AAATTTTCATC CTTTCTTTT TGATATTCAG AGCGATAAAA ATCCGTTTTT	19800
TGAAGTTTTT AAAGTTCCGA AAACCAAGG CATTCGCTT GATAAGTTTG ATGAGATTAT	19860
TGGTCGCTTC CAGTTTGGCG TTAGAATAGT GTAGTGAAG GCGGTTGATA ATCTTTTCTT	19920
TATCTTTGAG GAAGGTTTTA AAGACAGTCT GAAAAATAGG ATGAACCTGC TTAAGATTGT	19980
CCTCAATAAG TCCGAAAAAT TTCTCTGGTT CCTTATCTG GAAGTGAAA AGCAAGAGTT	20040
GATAGAGCTG ATAGTGCTGT TTCAAGCTTT CGAATAGCT CAAAGCTTG TTTAAAATCT	20100
CTTTATGGT TAAGTGCTA CGAAAAATAG GACGATAAAA TCGCTTATCA CTCAGTTTAC	20160
GGCTATCTTG TTGAATGAGT TTCCAGTAGC CTTGATAG	20199

(2) INFORMATION FOR SEQ ID NO: 7:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 19702 base pairs

(B) TYPE: nucleic acid

(C) STRANDEDNESS: double

(D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 7:

ACCCGATGTA TCAGCGGATA TTACTCTAT TTTCAAACG ATGTTATACC CACAATAAAA	60
GAAAAAGAC CCTAAGGTCT CTTTGCTTT TATTATTAAA CGCGTTCAAC TTTACCTGAT	120
TTCAAAGCAC GAGCTGAAG CCAACTTTT TTAGGTTTAC CATCGATAAG AACAGTAACT	180
TTTGAAGGT TTGGTTTAC GGCACGTTTT GTTTGGTTCA TCGCGTGTA ACGGTTGTTT	240
CCTGATACAG TCTTACGACC GTTAAAGTAA CATACTTTAG CCATTGTGTT TTCCCTCTAT	300
TAGATCTAAT ATAGCGGATG TGCTAGCACC ACATACCGTA CTATGTTATC ACATTTTCTT	360
GTTTTTTGCA AGGGAATTGG AAGATTTTTT ATTTGTGTCT TAAATCAGGT CTTCGCTGAC	420
ATTTCTGCTC TCCACATGCC ATCGTTGATT AACAGAACAC CAGAATTAAA ATTATGTGTA	480
TAAAAATCAT CTCTAACTGC AGCTAAGGGT ATAGCGGTCA AGTCCAATC CCACAGCTCA	540



TCTATCGATT TTCTTACAAC AATATCTGAA TCCAAATACA GTACACGAGA CTGCTTACA	600
TACTTTTGGAA TAAATACCT AAAAAAGCG CATATGAAAG TCCCTCAAAG GGGAGACGAT	660
AACCTTTCAG AATATTACTG TCAATCTAAA CATTACAAT CTCACATTC AAAGTCTCTA	720
GTCTTTTTC CATCAATTGG AACATTCTC GCGGAAGTC ATCATATAAA ACATAAACT	780
TAGATTATA ATGATGAACA CAAAGAGATT TATTGTTGT TTCAACTTTA TCCATATAAG	840
CATTATCTGC ACCTAAGACA ATCGCTTTT TCTCTCTTT CACTTTTAT CTCATTTCTT	900
TTTATTTCCA TCATATTATT CCAATCATAT GTTCCCATC ATATGTTCT ACGTAACCAT	960
TATTTTGCC TATTGCTCG TAAACCATTA CCAATGAGA TTTTAGATGA AGTCCCATTA	1020
CGGTTTACAA TTTTACATT ACGACACGGA GTTTTACAA TCGATTTCAT TTGCCAAAG	1080
TAGTTAGTGA GGCAGTTAGC TAGTTCGCCA AATAGCGACT AGCGTCCAAC AATTGGAAC	1140
TTTAGTTCCA ATTGTTGGTA CTGAGTCACA TCTCTCTCT TAACTCTAGC TCTGGATACT	1200
TGTCGCAAA CAGCGGAGG GCAAAGTCAT TTTCAAAGAG AAAGACTGGT TGGTCAAAC	1260
GGTCTTTGGC TAAGATATTG CGACTTGAGC ACATCCGTC ATCCAGTCC TCAGGCTTGA	1320
TTCAAAGAAC GGTCTTTTFA CCAATGGGT TCATAACTAC TTCGCAATTG TACTGGCCTT	1380
CCATGCGGTG TTTAAAGACT TCAAACCTGA GTTGACCTAC AGCGCCTAGC ATGTACTCAC	1440
CTGTTTGTA ATTCTTATAA AGCTGAACGG CTCCTTCTTG CACCAATTGC TCAATCCCT	1500
TGTGGAAGGA TTTTGCTTC ATAACTTCT TAGCAGAAAC TTTTCATGAA ATCTCAGGTG	1560
TAAAGTTGG CAGGGTTTCA AATTCAAAC TGTTTTTTC AACCGTCAAG GTATCCCCAA	1620
CTTGATAGT ACCGGTATCG TAAACCCCGA TAATATCACC TGCCACGGCA TTGGTCACAT	1680
TCTCACGACT CTCCGCCATA AACTGGGTAA CATTAGATAG TTTAGCCCC TTACAGTAC	1740
GAGGAGATT GACACTCTNG CCGGCTCAA ATTCCGCAAG TACGATAAG ACAAGGCCAA	1800
TACGCTCAGC GTGACGAGG TCCATGTTGG CTGGATTTT AAAGACAAAG CCTGAGAAAT	1860
CCTTGCTATA AGGATCCACA ATTTACCGT CTGTTTCTT GTGACATGT GGTCTGGAG	1920
CAAACTTGA GAAGTTTCA AGGAAGGTCT GCACAACAA GTTGTCTAGG GCTGAACCGA	1980
AAAAGACAG CGTCAATTCT CCAGCCAGAA TAGCTTCTCT TGAAACTCA TTCCCGGCTT	2040
CATTTAAAG CTCAATGTCA TCTTGACTT GCTCGTAGAA AGGATTGCTA CCAAGAGTT	2100
TGTCCCCGTC TTCTAGACTG GCAAAACGCT CATCCCCCT GTAAAGCTCT AAACGTTGCT	2160
TATAGAGGTC ATACAAGCCC TCAAGGCTT TCCCATCCC GATAGGCCAG TTCATAGGCT	2220
AGCTAGCAAT GCCCAAGATT TCTTCCAATT CTGCAAGAG ATCCAAAGGC TCACGACCGT	2280

		196	
CACGGTCCAG	CTTGTTTCATA	AAGGTAAGA	CTGGAATGCC
ACAAATTTCTT	GGTTTGAGCC	TGGATCCCTT	TGGCAGATC
CCACGGCCAT	CAAGGTACGA	TAGGTATCTT	CTGAGAATC
AGATAATTCAT	GCCTTGGCC	TGCTAGTCAA	ATTGCATAAC
CACGTTGCTT	CTCGATATCC	ATCCAGTCAG	ATTTAGCAAA
TTACCGTACC	AGCCTCACGA	ATCTCACCCC	CAAGTAGAG
TTTCCCGCG	GTCCGGGTGG	GAGATAATGG	CAAGGATAG
GAATAATTCAT	AAGTCTCTT	TCTTTGATTC	TCTATTTTTC
GATTTTACA	TTGATTTTA	CCATTCTTTT	CAACACTCCA
TTTTTCAAT	TCTATTCTT	TTCACTTCCC	CCTCCCTTAT
ATGAACAGA	CTAAAAATCA	TCATTTACAG	AAAGGATGCA
GAGGTAACAC	ACGTTGCCAA	TCTTTCAAAA	TTAAGATTCT
TTTGCGACCA	CCTTGTCTAA	GATTGTTGAC	ATGGTTGAAT
ACTGGTGTGG	CACCTACTAC	GACTATGGCT	GACCGCAAGA
GCGGAAGAG	GAATAGACCG	TGATCGCTTG	TTTAAAAACG
TATATCAAGG	TGCCAGCTAT	CCTAGACAAAT	GGAGGAGATG
AAACTATTGA	AGAGTTGCAC	AATCTCCTTG	TCTCTAAGGA
CCCAAGCAAC	ACTTGAATAT	ATCAAGTCTC	GTGAGGAAGC
TGCGTGAGGA	GCAAGCTCTT	GTTCAGGCTA	AAGCCATTGA
ACAATGTCTT	TTCAGGAATT	CCACTTGGCT	TTAAGGATAA
TCACAACTGC	TGCCCTAAAA	ATGCTCTACA	ACTATGAGCC
TTCCCAATGC	AAAAACCAAG	GGCATGATTG	TGTTGGAAA
CTATGGGTGG	TTCAGGTGAA	ACTTCAACAT	ACGGAGCAAC
GCAAGGTTCC	TGGTGGGTGA	TCAAGTGGTT	CTGCGCAGC
GCTTGTCAC	TGGTTCTGAT	ACTGGTGGTT	CCATCCGCCA
TGCTTGGTCT	CAAAACCAAC	TACGGAACAG	TTTCAGGTTT
GCTCATTAGA	CCAGATTGGA	CCTTTTGGCT	CTACTGTTAA
ACGCTATTGC	CAGCGAAGAT	GCTAAAGACT	CTACTTCTGC
TTACTTCAAA	AATCGGCCAA	GACATCAAGG	GTATCAAAAT
TAGCGCAAGG	AATTGATCCA	GAGGTTAAGG	AAACAATCTT

2340

2400

2460

2520

2580

2640

2700

2760

2820

2880

2940

3000

3060

3120

3180

3240

3300

3360

3420

3480

3540

3600

3660

3720

3780

3840

3900

3960

4020

4080

AAAAATGGG	TGCTATCGTC	GAAGAAGTCA	GCCTTCCTCA	CTCTAAATAC	GGTGTTCGG	4140
TTTATTACAT	CATCGCTTCA	TCAGAAGCTT	CATCAAACTT	GCAACGCTTC	GACGGTATCC	4200
GTACGGCTA	TGCGCGAGAA	GATGCAACCA	ACCTTGATGA	AATCTATGTA	AACAGCGAA	4260
GCCAAAGTTT	TGGTCAAGAG	GTAACAACCT	GTATCATGCT	GGGTACTTTC	AGTCTTTCAT	4320
CAGGTIACATA	TGATGCGTAC	TACAAAAAGG	CTGGTCAAGT	CCGTACCCCTC	ATCATTCAG	4380
ATTTGCGAAA	AGTCTTCGCG	GATTACGATT	TGATTTTGGG	TCCAACCTGCT	CCAAGTGTTG	4440
CCTATGACTT	GGATTCCTCTC	AACCATGACC	CAGTTGCCAT	GTACTTAGCC	GACCTATTGA	4500
CCATACCTGT	AACTTGCGCA	GGACTGCCTG	GAATTCGAT	TCTGCTGGA	TTCTCTCAAG	4560
GTCTACCTGT	CGGACTCCAA	TTGATTGGTC	CCAAGTACTC	TGAGGAAACC	ATTTACCAAG	4620
CTGCTGCTGC	TTTGAAGCA	ACAACAGACT	ACCACAACA	ACAACCCGTG	ATTTTGGAG	4680
GTGACIACATA	ATGAACCTTG	AAACAGTCAT	CGGACTTGAA	GTCCACGTAG	AGCTCAACAC	4740
CAATTCMAAA	ATCTTCTCAC	CTACTTCTGC	CCACTTTGGA	AATGACCAAA	ATGCCAACAC	4800
TAACGTGATT	GACTGTCTTT	TCCCAGGAGT	TCTACCAGTT	CTCAATAAAG	GGGTGTGTGA	4860
TGCGCGTATC	AAGGCTGCTC	TTGCCCTCAA	CATGGACATC	CACAAAAAGA	TGCACCTTGA	4920
CGCAAGAAC	TACTTCTATC	CTGATAACCC	CAAGCCCTAC	CAAAATTCCTC	AGTTTGATGA	4980
ACCAATCGGA	TATAATGGCT	GGATTGAAAT	CAAACTAGAA	GATCGTACGA	CCAAGAAAT	5040
CGGTATCGAA	CGTGCCACC	TAGAGGAAGA	CGCTGGTAAA	AACACCCATG	GTACAGATGG	5100
CTACTCTTAT	GTGACCTCA	ACCGCCAAGG	GGTCCCTTG	ATTGAGATTG	TATCTGAGGC	5160
AGATATCGGT	TCTCTGAAG	AAGCCTATGC	TTATCTGACA	GCCTCAAGG	AAGTTATCCA	5220
GTACGCTGGC	ATTTCTGACG	TTAAGATGGA	GGAGGTTGG	ATCGGTGTGG	ATGCCAACAT	5280
CTCCCTTCGT	CCTTATGGTC	AAGAGAAATT	CGGTACCAAG	ACTGAATTGA	AGAACCTCAA	5340
CTCCTTCTCA	AACGTTGCTA	AAGGCTTTGA	ATACGAAGTC	CAACGCCAGG	CTGAAATTCCT	5400
TGCGTCAAGT	GGTCAAAATC	GCCAAAGAAC	ACGCGCTAC	GATGAAGCGA	ATAAAGCAAC	5460
CATCTCATGT	CGGTCAAGG	AAGGGGCTGC	TGACTACCGC	TACTTCCCAG	AACCAGACCT	5520
ACCCCTCTTT	GAAATTTCTG	ACGAGTGGAT	TGAGGAAATG	CGGACTGAGT	TGCCAGAGTT	5580
TCCAAAGAAA	CGTCGTGCGC	GTTATGTATC	TGACCTTGOT	TTATCAGACT	ACGATGCTAG	5640
TCAGTTGACT	GCTAATAAAG	TCATCTCTGA	CTTCTTTGAA	AAGGCTGTGG	CCCTAGGTGG	5700
TGATGCCAAA	CAAGTCTCTA	ACTGGCTCCA	AGGGGAAGTC	GCTCAGTTCT	TGAATGCTGA	5760
AGGTAAACA	CTGGAACAAA	TCGAATTGAC	ACCAGAAAAAC	TTGGTTGAAA	TGATTCGCAT	5820

198			
CATCGAAGAC	GGTACTATT	CATCTAAGAT	TGCCAAGAAA
AAATGGCGTG	GGCGCGCTG	AATACGTGGA	AAAGCAGGT
AGCTATCTTG	ATCCCAATCA	TCCACCAAGT	CTTTGCCGAT
CTTCAAGTCA	GGCAAAAGTA	ACGCGACAA	GGCTTTACAG
AAAGGCCAAG	CCAACCCACA	AGTTGCCCTT	AACTACTTG
AAAGAAAAT	AGACAGAACA	AAACCAGCC	TAAGGTGGT
CCAATACTA	TTTTGGCTTT	ATTTCAGAG	TATTTATGG
TTTATTAAG	AGGTAAAAC	ATGATTGAAG	CAAGTACCT
AAACAGCTGA	CGGCAATTG	ATTGGCGTT	TGGAAGCTAG
GAAACAGAT	CATGCGTATG	AAATTGCGTG	ATGTCCGTAC
GCTACCGTCC	AGAGGAAAA	TTTGAACAAG	CTATTATCGA
TGTACAAAAT	GGATGACACA	GCATACTTCA	TGAATACAGA
TCCCTGTAGT	CAATGTTGAA	AACGAATTGC	TTTACATCCT
TCCAATTCTA	CGGAACGTAA	GTGATCGGTG	TCAACGTTCC
TTGCTGAAG	TCAACCATCT	ATCAAAAGTG	CTACTGTTAC
CGATGGAAAC	TGGACTTGTC	GTAACGTTT	CAGACTTCAT
TTATCAACAC	TGCAGAAGGA	ACTTACGTTT	TCCTGCTCTA
CTATGGGAAT	TGAAGAACAA	CTTGGCGAAA	TGGTTATGCG
TCATTTGCTAT	CGCTACTGCA	AAGGTAGAGG	GTGTTCACCT
CTGATACCTT	TTCAAAACTT	TCACTCGGCC	GTGGCATTTA
AACTCACAGC	AGATATCTAT	CTCTACCTTG	AGTACGGAGT
TTGCTATCCA	GAAGCTGTC	AAAGATGCCG	TCCGTAATAT
CTATCAATAT	TCACGTGCA	GGTATGTCC	CAGATAAAAC
ATCTAATTGA	CGAGGACTTC	CTCAATGACT	AGTCCACTAT
CGTAAATGCG	CTTTTCAMCG	TCTCATGAGC	CTTGAGTTTG
TGTGCTTTTG	CCTATACTCA	TGATCGTGAA	GATACGATG
ATAGACTCTG	TTTCTGTGTT	TCAAGCTAAA	AAGGAAGAAC
CATTTAAGAG	CAGGTGGAC	CAITGAACGC	TTAACGCTCG
TTGGAGTCT	TTGAAATCAC	TTCAATTGAC	ACTCTCTCAG
ATCGAGCTTG	CAAAGGACTT	CTCCGATCAA	AAATCTGCC

5880  
5940  
6000  
6060  
6120  
6180  
6240  
6300  
6360  
6420  
6480  
6540  
6600  
6660  
6720  
6780  
6840  
6900  
6960  
7020  
7080  
7140  
7200  
7260  
7320  
7380  
7440  
7500  
7560  
7620

AGCCAGTTTG	TAACAGAAGA	ACAATAAGGC	TCTTTGTCAA	CTGTAGTGGG	TTGAAAAAAA	7680
GCTAAGCTCG	AGAAAGGACA	AATTTCGTCC	TTTCTTTTIT	GATGPTTCAA	GGGATAAAAA	7740
TCCGTTTTTT	GAAGTTTICA	AAGTTTCGAA	AACCAAGGC	ATTGCGCTTG	ATAAGTTTGA	7800
TGAGATTATT	GGTCGCTTCC	AGTTTGGCAT	TAGAATAGTG	TAGTTGAAGG	GCCTTGACAA	7860
TCTTTTCTTT	ATCTTTGAGG	AAGTTTTTAA	AGACAGTCTG	AAAAATAGGA	TGAGCCTGCT	7920
TAAGATTGTC	CTCAATAAGT	CCGAAAAAAT	TCTCTGGTTC	CTTATTCTGG	AAGTGAAACA	7980
GCAAGAGCTG	ATAGAGCTGA	TAGTGGTGTT	TCAAGTCTTG	TGAATGGCTC	AAAAGCTTGT	8040
CTAAATCTC	TTTATTGGTT	AAGTGCAATC	GAAAGTAGG	ACGATAAAAT	CGCTTATCAC	8100
TCAGTCTACG	GCTATCCTGT	TGAATGAGTT	TCAGTAGCG	CTTGATATCC	TTGTATTTCAT	8160
GGGATTTCG	ATGAACTGA	TTTCATGATT	GGACACGCAC	ACGACTCATG	GCACGGCTAA	8220
GATGTTGTAC	AATGTGAAG	CGATCAAGAA	CGATTTTAGC	ATTCGGGAGT	GAAACAGTCT	8280
GGGAGACTGT	TTTCAGCCTGA	GCCTAGGAAT	TTGAAAGCGA	AGCTGTTTAG	CCAAGTCATA	8340
GTAAGGGCTA	AACATATCCA	TAGTAATAAT	TTTGACGCGA	CATCGACAA	CTCTATCGTA	8400
GCGAAGAAAG	TGATTTCCGA	TGATAGCTTG	TGTTCTACCC	TCAAGAACAG	TGATGATATT	8460
GAGATTGTTA	AAATCTTGCG	CAATGAAGCT	CATCTTTCCC	TTTGTAAAAG	CATACTATTC	8520
CCAAGACATA	ATCTCAGGAA	GACAAGAAAA	ATCATGTTTA	AAGTGAAAT	CATTGAGCTT	8580
ACGAATAACA	GTGAAAGTTG	AGATGGAAG	CTGATGGGCA	ATATCAGTCA	TAGAAATCTT	8640
TTCAATCAAC	TTTTCAGCAA	TCPTTTGGTT	GATGATACGA	GGGATTTGGT	GATTTTCTT	8700
GACGATAGAA	GTTTCAGCGA	CCATCATTTT	TGAACAGTGA	TAGCACTTGA	ATCGACGCTT	8760
TCTAAGGAGA	ATTCTAGTAG	GCATACCAAT	CGTTTCAAGA	TAAGGAATTT	TAGAAGTTTT	8820
TTGAAAGTCA	TATTTCTTCA	ATTGTTTCCC	GCACTCAGGG	CAAGATGGGG	COTCGTAGTC	8880
CAGTTTGGCG	ATGATTTCCCT	TGTGTGTATC	CTTATTGATG	ATGCTAAAAA	TCTGGATATT	8940
AGGGTCTTTA	ATGCTAGTGA	ATTTTGTGAT	AAATGTAAAT	TGTTCCAAT	GAATCTTTCT	9000
AATGAGTTGT	TTTGTGCTTT	TTCATTATAG	GTATATATGG	ACTTTTTTTC	TACAATAAAA	9060
TAGGCTCCAT	AATATCTATA	GGGGATTTAC	CCACTACAAA	TATTTATAGAG	CCAACAATAA	9120
AAAGAAAAAG	TGTTTGAATG	ATATCAAAACA	CTTTTTCCTT	TGCTTCCAC	TATCTAAAAA	9180
AATGATAATA	GATATAAATTG	TAAACAAAAA	TCCAGATAGG	TTTTGCATGA	TTGAGAAAGT	9240
TAAAAAAACT	ATGGCAGAGA	ATCGTTAATC	TCCAGATTGC	GGTAGAAGCA	TAAACAAGGG	9300
CAAAAAAGAA	ACCAATCAGA	CTATAATATA	ATAAATAAT	TGGATCTCTG	TGAGATAGTA	9360

	200	
TCAAATGGCT AATCCCAAAG ATGATAGCAG ATAGGATAAC ATCCAAATAG TACTTTGGACT	9420	
AGGGAJAGAA GGTATTTCATA AAATACCCCTC TATCAAGAGT CTCTCCAAAA ACAGGACCGA	9480	
TGATTACAGG CAGGACAAAA GATAAGATAG TCGATAAJAA GGTGGGTGT CCATTTGAAA	9540	
AAAGCACGGT AAAATACTCA TCATGAATAT TCCTATGAAT AATCAAATGA GCATAGCGTG	9600	
CCCAAAATTT ACCGAGAATC TGAATAACCA CATAGATTGC AAATAAGTAG AAGACAAATG	9660	
ACCAGTTCCA GCTCTTTTTC TCAAGAGATA AGAGCATCTT TTCTTTTTF AACCTCCAAA	9720	
TTAATAGAGG GAAACTTCCC ACTAATCCCA TTGTTAAAT AAGAGAATAG ACATCAGCTC	9780	
CTAACCTTAA ATGATCGTC ACATACAATC CAATTTTTC TGGTAAATAG GTAGATAGTA	9840	
AAATATAAG CAAAAATATT CCAAATTTGC TTAGTTTTTT TGTGTTTCTC ATCGTACTTT	9900	
TTTGAAGAT TACCTGCTC GGAAGCGTA CTCCAGCA TCTATATAAG AATTAAGTGC	9960	
CCCTTGCTC ATATAGGGAG CAATTTCTCT ATAATAAAC CATCTACTAT ATCATCTTC	10020	
CCAAACAGCA AGACCACCTG AAGTTTGCTC CAAGCTCTCA GTTGAAGAA CTCTAAATGT	10080	
ATTTGTACTT GTCATTTGCA GTACCTTCTT AAAATAGATT GTTGTAGGCT CACATTTATA	10140	
GTATATTTCT TTTTTTGTCT ATTTTATAGC CCATCTCCTC AACTGGCAAT TTTTCGACCT	10200	
GAATTACATT TTTCCATAAA AAATGAGACC TTCTTAGTCT CATTTAGTCA TTCTTAGTAT	10260	
TTTCTAAATC GTTGATAGCG TTCTTCCAGC AACTCTTCTA GCGTTTITG TGAAAGTCTA	10320	
GCCAGCTCG TTTGGAGTTC TTTTTTGACA CTCTTAATCA GTTCTTTACT AGAAGTCTCT	10380	
ATTTTCAGAA TCACCTTATC CACCACGTCC ATTTCTAACA GTTCATCGGA AGTGATTTTC	10440	
ATGAGTTCTG CTGCTTCCAT AGCGCGAGTA CGTCCCTGCC ATAAAAAGGA AGCAAGSCT	10500	
TCTGGACTGA GAATGGCATA GATAGAAATT TCCAGCATCC AGACACGGTC CGCGACAGT	10560	
AGAGCCAGAG CCCCGCCTGA ACCACCTTCA CGATAATAA TGGGATAAT AGGAACCTTC	10620	
AGGTCACTCA TTTCCATGAG ATTGCGAGCG ATAGCTTCCC CTTGACCAG TTCTTCCGCT	10680	
CCGACACCAG GATTAAGCAC TGCTGTATTG ATAAAGGTCA CAACTGGACG GCCAAATTTT	10740	
TCAGCCTGTT TCATCAACCG CAGTGCCCTT CGGTAGCCTT CTGGATGTGG TTGGCCAJAA	10800	
TTCCGTTTGA GGTGTCTTTG CAATCTCTTG CCTTTTTTGA TACCAACCA TGTTPACAGT	10860	
TGCTCTCCAA GCCAACCAAT ACCACCAACA ACTGCACCAT CATCACGAAA AGAAGGCTCA	10920	
CCATGTAAAT GATATAATTC ATCAAAAATG CCTGTGCAAA AGTCAAGGT TGTCAAGCTA	10980	
CTCTGCTCAC GCGCTTCTCT GACTATTTTT GCAATATTCA TCTAGGACTC CCTCCATGCA	11040	
ATCTGACTAG GCTAGCAATC GTATCTGGTA AGTCTCTTCT TTTGACAAAT GCATCCACAA	11100	
AGCCATGTTT TAATAGGAAT TCTGCCCTTT GGAATCCTC AGGCAAGCTT TCACGAACCG	11160	

TATTTTCAAT CACACGACG CACGAAAAC CAACCAAGCT CTGTGGTTCA GCCAGAATGA	11220
TATCGCCTTC CATAGCGAAA GAAGCTGTCA CACCACAGT CGTTGGATCT GTCAAATGK	11280
TCAGTAAAA GAGACGACA TTGGAATGGC GTTAAACCG CGCAGAGATC TTAGCCATCT	11340
GCATGAGACT CATGATTCTT TCCTGCATAC GGGCTCCACC AGAGGCTGTG AATAGGACAA	11400
CTGGCAATTT TTGACAGTC GCATCTCAA ACAACAGT GATTTTTCCT CCTACAACCG	11460
TACCCATAGA AGCCATGATA AAGTTAGAAT CCATTAATCCC AAGAGCCACA GTCTGACCTT	11520
TAATAAGAGC AGTTCTCTGC ACAACGGCTT CATGCAAGCC TGTTTTTCCT CGCATAGATG	11580
CCAGTTTCTT TTGTTAACCA GGGAAATGCA AGGGATCCTT GCCTTCAATC CCTGTAAACA	11640
ATTCTTTTGA GGTTCGCATA TCAATCGTCA AAGCAAGCG TTCTTGGGCA GAAATACGAA	11700
AGGTATAGCT ACAGTCCGGA CAGATACGTT CACTTCCCAG ATCCTTCTGA TAGATGGTAT	11760
GCTTACAGCC TCGACACTGG GAAAATAATT CATCTGGAAC CTCGTGGCTTA GCTTGAGGTT	11820
TTTCCCTAAC CGAACGATTG GGATTGATTC GAATATACTT ATCTTTTFTA CTAATPAGAG	11880
CCATTGATTC CCTTTTGGG TTTAACTCTT TAAAGTCATT TTATTCTTTT TCTTGATATT	11940
TAGGTAAAGAA GGTTCCTATC AAGAAGGAAG TATCATATCT CCCAGCAATC ACATTGCGAT	12000
CTGAATGAG GTCAAGCTGG AATCTGCAAT TGGTCTGCAC TCCTTCAATT TCTAATTCAT	12060
AGAGGGCAGC TTGCATTTTC ATCAAGGCGT CAAACGATT TTCCGCGTGT ACTATGATTT	12120
TGGCAATCAT ACTATCATAA TAAGCGGAA TGGTATAACC TGGATAAATC GCTGAATCCA	12180
CGCGCAAGCC AACTCCACCA CTGGGCAGAT AGAGATTAGT AATCTTACCT GGACTTGGAG	12240
CAAGTTAAA GGTCTGGTTT TCTGCATTGA TACGACACTC GATGGCATGA CCGGTAGGA	12300
CAATATCTTC TTGCTTAACA GACAAAGGCT GACCTGCCGC AATGCAAAATC TGTTCCTTAA	12360
CGATATCAAC ACCTGAACA AACTCTGTTA CTGGATGTC TACCTGAACA CGAGTATTCA	12420
TCTCCATGAA ATAGAAATG CTACTTGTCT CATCAAGAG AATTCATATG GTTCTGTCAT	12480
TCTCATAGCC AACAACTCT GCCGCTCGAA CAGCAGCAGC ACCTATTTC TGACGCAGCG	12540
TTTTTCCGAT TGCAATCGAG GGACTTTCTT CCAAAACCTT TTGGTTATTC CTTTGAAGAG	12600
AACAATCCCG TTCAACCAAG TGAATCACAT GTCCATGCTC ATCACCATTG ATTTGAACCT	12660
CAATGTGCCG AGCTGGATAG ATAAACCGTT CTATGTACAT GGCACCATTC CCATAATTGG	12720
CCTTGGCCTC ACTAGAGGCA GTTTCBAAGG CAGAAACGAG GTCATCTGCT TTTTCAACCT	12780
TACGAATCCC TTTTACCACT CCACCTGCTG AAGCCTTGG CATTAACAGA TAGCCAAATT	12840
TTTCAACAAC ATCAAAAGCT TCTTCAGAT TATGCACCTC TCCAATCGAA CCTGGTATAA	12900

202

CAGGCACACC	TGCTTTAATC	ATCTGAGCAC	GCQCTTGAT	CTTATCCCCC	ATCATATCCA	12960
TAACATGACC	AGATGGACCG	ATAAATCTGA	TACCTACTTC	TTCCACACATG	GTCCGCAAAAT	13020
TGGAATTTTC	ACTGAGAAAT	CCAAAACCAAG	GGTGAATAGC	TTCTTGCTTCA	GTCAAGACTG	13080
CAGCTGATAG	AACTGCATTA	ATATTGAGAT	AAGACTCTGT	TGCCTTGCCA	GGACCAATAC	13140
AAACTGCTTC	ATCTGCCAAA	AGCGTATGAA	GAGCTTCCTT	ATCAGCAGTT	GAATAAACCG	13200
CTACCGTCGC	AATCCCAAT	TCAOCTGCGC	CACGATAAAT	ACGAACCGCA	ATTCACCCAC	13260
GATTGGCAAT	TAAATTTTTT	CGAAACATGG	AGAACCTCCT	TAGTTCCCAA	TTGCAAAAGT	13320
AAGGGTACCA	CTGGGTGCAA	GCTTGCCATC	CACCTCAGCC	TTTGCTTCAA	CCACAGCTAT	13380
GGTGCCACGA	CGTTTTACAA	AAGTCGCTGT	CATAACCAAT	TGGTCGCTG	GTACAACTTG	13440
CTTCTTGAAAC	TTAACTTTGT	CCATACCAGC	GTAAAAGACC	AGTTTTCCTT	TATTTTCAGG	13500
TTTTGATAAC	TCCAACACAC	CGGCAAGCTT	CGCCAAGGCT	TCCATAATCA	CAACACCTGG	13560
CATAACTGGG	TATTGAGGAA	AGTGGCCGTT	AAAGAAAGGC	TCGTTGATGG	TCACATTTTT	13620
GATAGACAAC	ATGGTATCCT	CGCTCACTTC	CAAGACACGG	TCCACTAGAA	GCATAGGATA	13680
ACGGTGGGGA	AGAGCTTCTT	TGATTCTTTG	AATATCGATC	ATTTGATACG	TACCAATCCT	13740
TTACCAAACT	CAACCATTTT	TTCTTTAGAG	ACGAGAATTT	CGTTACCAC	ACCATCTTAA	13800
GGAGCTGGGA	TTTCAATCAT	GACTTTCATG	GCTTCGATAA	TTACCAATGT	TTGACCTTTT	13860
TTGACACTAT	CACCAACTGT	AACGAAGGCA	GGTTTATCTG	GTCCAGCAGC	CAAGTAAACC	13920
ACTCCAACAA	GTGGACTCTC	TACAAGATTT	CCCTCAGTAG	CCACACTTGC	TTCAGCTGGA	13980
GCTGGAACCT	CTTCTGCTAC	AGTCTCTGCT	GGAGCAGATG	TAGSAGCTAC	TGGACTCGGT	14040
GTGTGTAGAA	CGGGTGTCTG	AGCGACTTGA	GTTCGCAACT	CAGGCACAGG	TCTTGCTTCA	14100
TTCTTGCTAA	ACTGCAACTC	ATCCGTGCCA	TTTTTATAAG	AAAATCTCTC	CAAACCTGAC	14160
TGGTCAAAAT	GAGTCATCAA	GTCTTTAATA	TGCTTTAAAT	TCATACTTAT	CTATTCTCCC	14220
AACGTTTGAA	AGCAAGAACT	GCATTGTGGC	CTCCAAAACC	AAAAGTATTT	GAAATAGCGT	14280
ATGGAATTTTC	TTTCTCCAAG	CCTTGTCCAT	AAACGACATT	AGCTTGATTA	TAATCTGATA	14340
CTTCACITGT	CCCAGCTGTC	ATTGGTACAA	AGTTATGACG	CATAGCTTCG	ATGGTGACGA	14400
TAGCTTCTAC	TGCACCCGCA	GCCCCCAGCA	AATGTCCTGT	AAAAGACTTG	GTGTGATGATA	14460
CAGGTACTTTC	CTTACCAAGA	ACAGCTACGA	TAGCACCACT	TTCTCTTTTT	TCATTGGCAG	14520
GAGTTGACGT	TCCGTGAGCA	TTGACATAGG	CTACTTGCTC	TGGAGAATTC	TCAGCTTCTT	14580
CCAGGCTAG	TTTGATGGCC	TTGATAGCTC	CCTGACCTTC	TGATATGGGA	GAAGTCATGT	14640
GGTAGGCATC	ACAAGTATTT	CGTAACCAA	CCACTTCAGC	CAGGATAGTA	GCTCCAGGTT	14700



TTTCAGCGTG TTCAAGACTT TCTAGAACCA ACATCCCTGA ACCTTCACCC ATAACAAACC	14760
CATTCGGATC CTTATCAAAT GGGATCGAAG CACGAGTTGG ATCCTCTGTA GTAGAGAGAG	14820
CTCTTAAGGC TTGGAACCA GCGATGGCAA AAGGTCTGAT AGAAGCTTCT GTTCCTCCCA	14880
CCAACATCAC ATCTTGGAAA CCAACTTAA TGGAGCGGAA GGCATCCCCA ATCGCATCAT	14940
TTGATGAGA GCAGGCGATG TTGATAGATT TACAACACC GTTTGCACCA AAACGCATGG	15000
CTACATTCOC AGAAGCCATA TTTGGTAAAG CTTTGGGAAG AGTCATTGGT TTGACACGTT	15060
TGGGTCCCTT TTCATGAAGG CGAAGTACCT GATCTTCAAT TTCTTGAAT CCACCAATAC	15120
CAGATGCAAC GATAACACCA AAACGATCCC TATTAGAGC CTCTACATCA AGATTGGCAT	15180
GATTTCAGCG CTCTTGGGCT GCATACGAGG CATATAAAGA ATAGTTATCA AAACGGTTGG	15240
TATCTTTTTT TACAAAGTAT TTATCGAAGC GAAAACTTGG GATTCTGCC GCATTATGCA	15300
CATCAAGTGC ACTATGATCA AATTTTGTAA TGCCACCAAT GCCGATTTTC CCAATTGCTA	15360
AACTATTCCA AAATCTCTCT GGTGTATTTC CGATGGAGA TGTTACTCCA TAACCTGTTA	15420
CCACTACTCG ATTTAGTTTC ATCTTTTCA CCTCTAGCTT TCGCTACATA CTTAAGCCAC	15480
CATCAATGGC AACCACCTGT CCACTTAGAT AATCTTGGCC TGCTAAAAAT ACTGTCAAT	15540
CTGCAACCTG CTCTGCCCTG CCAATTTCTT TCATCGGAAT CTGAGCTAGT GTAGCTTCTT	15600
TAATCTTATC TGACAGGATA GCGGTCAAT CAGACTCAAT CATTCCTGGA GCAATCACAT	15660
TGACTCGTAT ATTCOGACTA GCGACCTGCG GTGCCACAGA CTTGGTAAAG CCAATCAAGC	15720
CAGCCTTAGA AGCAGCATAA TTAGCTTGAC CAATATTCCC CATCAAAACCA ACACACTAG	15780
ACATATTAAAT GATAGCACCT TCTCTGGCTT TCATCATCGG TTTCAGAGCT GATTGTGTCA	15840
TATTAAAGGC ACCAGTCAGA TTGACCTTGA GCATTTTTC AAAATCTGCT TCTGTCACTC	15900
TGAGCATAAG AGTATCTTGG GTAATCCCTG CATTTGTGAC CAAAACATCT ACTGAACCA	15960
GTTCGCAAT AGCTTGATCA ATCATACGCT TAGCGTCTGC AAAATCTGAT ACATCTCCTG	16020
AAATGGGAAC CACCTTGATA CCAATAGTTG AAAACTCAAG GAGCAATTCT TCTGAGATTG	16080
CCCCACGACT GTTTAAGACA ATGTTGGCTC CTGCTTGAGC AAACCTGTGG GCGATGGCAA	16140
GACCAATTCC ACGACTCGAA CCTGTAATAA AGATATTTTT ATGTTCTAGT TTCAATTTTT	16200
TCCTTTCAAA ACTTCTACTT ATTTTAGTCT ATTTTCTTAA AAGTGCTACT AAACCTGGCT	16260
GATCTTCAC ATGAGCTAAG TGACAGGTTT GATCAATTTT TTTAACAAAA CTTGACAAAG	16320
CTTTCOCGG TCCAATCTCG ATAAAGTTGC TTATGCTGCT TTCTTGCAAT ACCCAATAC	16380
TTTCATAGAA ACGAACGGGT TCCCTGACCT GACGCGTCAA GAGCTGAGCA ATGTCCTCTT	16440

TTTGCATCAC AGCAGCTTCT GTATTGCCGA CTAGGGGACA AGTAAATCT GAAAACTTA	16500
CCTGAGCTAG AGTTTCAGCT AGTTTCTGGC TAGCAGGTTC AAGGAGAGCG GTGTGAAAGG	16550
GACCTGACAC CTTAGAGGGA ATCAAGCGTT TGGCACCTGC TTCTTGCAAA AGTTCAACCG	16620
CTCGATCAAC TGCAACCACT TCTCCAGCAA TGACGATTG TGCAGGTGTG TTATAGTTGG	16680
CTGGAGTAAC CACTCCAAGT TCAGAAGCTT TTTGACAGGC TTCTTCAATG ACCTCTACTG	16740
GCGTATTGAG AACTGCTACC ATCTTGCCAG AGTCAGCAGG AGCCGCTTCT TCCATATAGT	16800
CTCCACGCTT AGCTACCAAG GCAACCGCAT CTTCAAAATC CAAGGGGCCA CTTGCCACCA	16860
AGGCAGAGTA TTCTCCAAGA GACAAACGAG CAACCATATC AGGCTGATAG CCCTTTTCTT	16920
GCAATAAAGC GTAGATAGCA ACCGAAGTCG CTAGAATGCG TGGTTGCGTA TAGCGGGTCT	16980
GATTGAGTTT GTCTTCTTCC GTATCGATGA GATAACGCAA ATCATAACCG AGCACCTGGC	17040
TGCGTCGATC AATCGTTTCT TTAACAATCG GATACTGATC ATAGAAATCC CGTCCCATCC	17100
CTAGATACCTG GGCACCTTGA CCAGCAATAA AAAAGGCTGT TTATGTCATT TCTTACAAC	17160
CCTGTCCAGC GAGAGGCTTC TTCTTGAATT TTCTTAGCGG CTCCGTAATA CAAATCTTTT	17220
AGGATTCTCT CAGCTGTGTC TTCTTTAGAA ACAAGCCCTG CGATTTGACC TGCCATAACA	17280
GAGCCACCAT CCACATCACC GTGAACAACT GCTTTGGCTA GAGCACCTGC TCCCATTTGT	17340
TCAAAGATT TTAATCAGG ATCTTCTTGC TTAAAGGCAT CTTTTCAGC CAGTTCAAAA	17400
TTCTAGTCA ACTGATTTT AATAGCACGA ACAGCATGAC CAAAGTGCTG AGCTGAAATC	17460
GTAGTATCAA TATCCCTTGC TTTTAAAAAT TTCTCTTGT AGTTTGGATG GGCATTGCAC	17520
TCTTTTGCAA CTACAAACCG TGTCCCCACC TGTACAGCCT CTGCACCTAG CATAAAGCCA	17580
GCCGACGACG CTTCAACATC CGCAATTCTC CTGCGAGCAA TAACAGGAAT AGATATAGCT	17640
GTGGCTACCT GTCGACCAA GGTCTAGGTT GTTAATTTAC CGATATGCCC CCCAGCTTCC	17700
ATTCTTCTG CAATAACAGC GTCTGCACCG ATTTTTTCCA TCGGTTTACG TAAAGCGACA	17760
CTAGGAACAA CAGGAATAAC GATTATCCCA GCTTCATGGA AACGTTCAT ATACTTGCTT	17820
GGATTTCCTG CTCTGTGTGT GACAACTTTA ACACCTTCTT CAATAACGAG ATCCACGATG	17880
TCTTCCACAA AGGGAGATAA GACGATGATG TTGACCCCAA AGGGTTTATC AGTCAATGAT	17940
TTGATTATTAT CAATATTGGC CTGACAACTT TCTTTGGGG CATTTCCCCC ACCGATAATT	18000
CCTAATCCTC CAGCCTTGA AACAGCCCTC GCCAAATCAC CATCAGCAAC CCAGGCCATC	18060
CCTCCTTGGA AAATAGGATA ATCAATCTTC AATAATTCTG TAATACGCGT TTTCATAGTG	18120
CCTCCAACCT TCCTTGCTTA CGTAATAGTT CGATTTCACC ATAAATTGAC AGTCAAACTA	18180
TTACCTAAAC AAGAGGAGT GGGTTTCTCC CTACTCCTTC TACTAATATT CTGCTAATT	18240

205

TGCTTGCTCT TCAACGTAAG CAACCAAGTC ACCAACTGTT TTCAAGTCAT TTCTGCTTC 18309  
 GATTGGGATA TCAAAAGCAT CTTCGATTTC TGAGATTACT TGGAAACAAG CCAATGAATC 18360  
 TGGCTCCAAA TCATCAAAAG TTGATTCAAG TGTTACTTCT GATGCGTCTT TTCCAAGTTC 18420  
 TTCAACGATA ATTTCTTGTA CTTTTTCAAA TACTGCCATG ATGAGACTCC TTTAAATTA 18480  
 ATAGTTTTTT TATAACAATG TGTTCACCAC ATGATTACCT AAATTTGAAG AATGAGCGTG 18540  
 CCCAGGTCA AGCCTCCACC GAAGCCTGAT AGAAGAACAG TCTGGCTACC ATCTAAAGGG 18600  
 ATGAGACCTT GTTCTACACA CTCTGAAAGT AAAATCGGGA TACTGGCTGC ACTGGTATTG 18660  
 CCATATCCCA TCATATTGGC TGGAGTTTG GCTCGGTCAA CACCAATTTT TCTAGCCATC 18720  
 TTATCCAAAA TACGGTCATT GGCTTGATGA AGTAGCAGAT AATCCAGTTC TGTACCTCT 18780  
 ATAGGAGATT CATCAATAGT CTGCTTGATA GACTTGGCTA CATCTCGAAT GGCAAAATCA 18840  
 AAGACTGTGC GTCCATCCAT CTTCAAAAAC GAATCTGCAC TTTCTTGATC TGAATATGGA 18900  
 GAATGTAAAC CTGAATGCCC ATAAGTTAAA CACTCGCTGC GACTTCCATC GCTATTGAGA 18960  
 CTCTCAGCTA AGAATGCTC TTGCTCGCTA GCTTCTAACA AGACACCACC AGCACCATCT 19020  
 CCAACACCA CAGCTGTGTA TCGATCCGAC CAATCGACTG CTTTAGAGAG GGTTCACCTA 19080  
 CCAATCACCA AGCCTTTTGG AAAGCGACCA GAAGCGATAA ACTTTTCAGC AGTTGAAAGA 19140  
 GCAAAATCAA ATCCACTGCA AGCCGCGGTT AAGTCAAAAG CAAGGCTTT ATTAGCACCA 19200  
 ATATTAGCTT GAACACGAGC AGCTGTAGAG GGCATCATCG AATCTGGAGT AATGGTAGCT 19260  
 AGGATGATAA AATCCAGTTC TTCTCCTGTT ATTCAGCTT TTGCCATCAG TTTCTTAGCA 19320  
 ACCCTCTGAG CCAATCACT GGTAGATTCT GTTCTTGAAA TATGCTTTG TCGTATTCCC 19380  
 GTTCGACTTG AAATCCACTC ATCATTTGTA TCCATAATCT GAGCCAAAGC GTGATTGTGA 19440  
 ACCACTTGCT CTGGCACATA ATGAGCAACC TGACTTATTT TTGCAAAAGC CATTATTCTA 19500  
 AATCCTCCAA AAATTGGTAA AGATTAGTCA AACCTTTACC CATGACAGCA ATTTCTTCCT 19560  
 CGCTCATGCC ATCAATAATT TTTTCTACCA TGGCCTTGTG GAAGCGTTTA TGCAGTCTAT 19620  
 GAATCAAGCG ACCCTTCTTT GTCAAAATGCA GATGCACCAC ACGACGATCC TGTCTGACC 19680  
 GAACCTGCTC AATGTAGCCC GG 19702

(2) INFORMATION FOR SEQ ID NO: 8:

(1) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 6211 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 8:

GAAAAATTCC TCTCTTCTCT TGAAAAAATT TGAAAAAATG GTATGATAGT AACAAAGTTAT	60
TTTTAAGCAGG AAAGAAAGGG GAATAATGGA GAAAAATCAGT TTAGAATCTC CTAAGACGGG	120
GTCGGACCTA GTTTTGGAAA CACTTCGNGA TTTAGGAGCT GATACCATCT TTGGTTATCC	180
TGGTGGTGGC GTTTTGCCCT TTTATGATGC GATATATAAT TTTAAAGCCA TTGCCACAT	240
TCTAGGCGGC CATGAGCAAG GTTGTTTGCA TGAAGCTGAA GGTATGCGCA AATCAACTGG	300
AAAGTTGGGT GTTGGCGTGC TCACCTAGTG ACCAGAGCA ACAAAATGCC TTACAGGGAT	360
TGGGATGCC ATGAGCGATA GCGTTCCCTT TTTGGTCTTT ACAGGTCAGG TGGCGCAGC	420
AGGGATTGGG AAGGATGCGT TTCAGGAGGC AGACATCGTG GGAATTACCA TGCCAATCAC	480
TAAGTACAAAT TACCAAGTTC GTGAGACAGC TGATATTCCG CGTATCATTA CGGAAGCTGT	540
CCATATCCCA ACTACAGGCC CTCAGAGGCC AGTGTAAAT GACCTACCAA AAGACATATC	600
TGCTTTAGAA ACAGACTTCA TTTATTCAAC AGAAGTGAAT TTACCAAGTT ATCAGCCGAC	660
TCTTGAGCCG AATGATATGC AAATCAAGAA AATCTTGAG CAATTGTCCA AGGCTAAAAA	720
GCCAGTCTTG TTAGCTGGTG GTGGAATTAG TTATGCTGAG GCTGCTACGG AACTAAATGA	780
ATTTGCAAGAA CGCTATCAAA TTCCAGTGGT AACCACTCTT TTGGGACAAAG GAACGATTGC	840
AACGAGTCCAC CACTCTTTTC TTGGAATGGG AGGCATGCAC GGTCAATTTC CAGCAAAATAT	900
TGCTATGACG GAAGCGGACT TTATGATTAG TATTGGTTCT CGTTTCGATG ACCGTTTGAC	960
GGGGAATCCT AAGACTTTTC CTAAGAAATGC TAAAGTTGCC CACATTGATA TTGACCCAGC	1020
TGAGATTGGC AAGATTATCA GTGCAGACAT TCCTGTAGTT GGAGATGCTA AGAAGGCCCT	1080
GCAAAATGTTG CTAGCAGAAC CAACAGTTCA CAACAACACT GAAAAGTGGA TTGAGAAAGT	1140
CACATAAGAC AAGAATCGTG TTCTGTTCTA TGATAAGAAA GAGCGTGTGG TTCAACCGCA	1200
AGCAGTTATT GAACGAATTG GTGAATTGAC GAATGAGAT GCCATTGTGG TAACAGACCT	1260
TGGTCAACAC CAAATGTGGA CAGCTCAGTA TTATCCCTAC CAAATGAAC GTCACTTAGT	1320
GACTTCAGGT GGTTTGGGAA CAATGGGCTT TGGAAATCCA GCAGCAATCG GTCTAAAAAT	1380
GCTTAACACA GATAAGGAAG TAGCTCTGTT TGTGGGGAT GGTGGTTTCC AAATGACCAA	1440
CCAAGAGTTG GCTATTTTGA ATATTTACAA GGTGCCAAAT AAGGTGGTTA TGCTGAACAA	1500
TCATTCACCT GGAATGGCTC GCCAGTGGCA GGAATCCTTC TATGAAGGCA GAACATCAGA	1560
GTGCGCTTTC GATACCCCTC CTGATTTCGA ATTGATGGCG CAGGCTTATG GTATTAAGAA	1620
CTATAAGTTT GACAAATCTG AGACCTTGGC TCAAGACCTT GAAATCATCA CTGAGGATGT	1680

TCCATATGCTA APTGAGGTAG ATATTTCTCG TAAGGAACAG GTGTACCAA TGGTACCGGC 1740  
 TGTFAAGAGT AATCATGAGA TGTTCGGGGT GCAGTTCCAT GCTFAGAATG TTAACAGCAA 1800  
 AACTACAAA TCGTTCAGGA GTCTCAATC GCTTTACAGS TGTCTATCT CGTCGTCAGG 1860  
 TPAAPATTGA AAGCATCTCT GTTGGAGCAA CAGAGATCC GAATGTATCG GGTATCACTA 1920  
 TTATTATTGA TGTTCCTTCT CATGATGAAG TGGAGCAAT CATCAACAG CTCATCGTC 1980  
 AGATTGATGT GATTGCGATT CGAGATATTA CAGACAAGCC TCATTGGAG CGCGAGGTGA 2040  
 TTTTGGTTAA GATGTCAGCG CCAGCTGAGA AGAGAGCTGA GATTTTAGCG ATTATTCAAC 2100  
 CTTTCCGTGC AACAGTAGTA GACGTAGCGC CAAGCTCGAT TACCATTAG ATGACGGAA 2160  
 ATGCAGAJJA GAGCGAAGCC CTATTGCGAG TCATTGCGCC ATACGGTATT CGCAATATTG 2220  
 CTCGAACGGG TGCAACTGGA TTACCCGCG ATTAATAATC CAACTAAAT TTATTAAACC 2280  
 AGCTTAAAG CCAATAAATA ATAGAAAGA GAGAAAGCT ATGACAGTTC AATGGAAATA 2340  
 TGAATAAGAT GTTAAAGTAG CAGCACTTGA CGGTAAAAA ATCGCGTTA TCGTATTAG 2400  
 TTCACAAGG CATGCGCATG CTCAAAACCT GCGTGAATCA GGTGCTGACG TTATTATCGG 2460  
 TGTACGTCCA GGTAAATCTT TTGATAAAGC AAGAAGATG GGATTTGATA CTACACAGT 2520  
 AGCAGAGCT ACTAAGTTGG CTGATGTTAT CATGATCTTG CGCCAGACG AATTCACAA 2580  
 AGAATTGTAC GAAGCAGAAA TCGCTCCAAA CTTGGAAGCT GGAACGCAG TTGGATTGTC 2640  
 CCATGGTTTC AACATCCACT TTGAATTTAT CAAAGTCTCT GCGGATGTAG ATGCTTTCAT 2700  
 GTGTGCTCCT AAGGACCAG GACACTTGGT ACCTCGTACT TACGAAGAAG GATTGGGTGT 2760  
 TCCAGCTCTT TATGCAATAT ACCAAGATGC AACAGGAAT GCTAAAAACA TTGCTATGGA 2820  
 CTGGTGTAAA GGTGTGGAG CGCGCTGTGT AGGTCTTCTT GAAACAACCT ACAAGAAGA 2880  
 AACTGAAGAA GATTGTGTTG GTGAACAAGC TGTACTTTGT GGTGGTTTGA CTGCCCTTAT 2940  
 CGAAGCAGGT TTGAACTCT TGACAGAAAC AGGTACGCT CCAGAAATGG CTTACTTTGA 3000  
 AGTCTTTCAC GAAATGAAT TGTATGTTGA CTTGATCTAC GAAGTGGAT TCAAGAAAT 3060  
 GGTCAATCT ATTTCAACA CTGCTGAATA CGGTGACTAT GTATCAGGTC CAAGTGAAT 3120  
 CACTGAACAA GTTAAAGAA ATATGAAGGC TGTCTTGGCA GACATCCAAA ATGGTAAAT 3180  
 TGCAAATGAC TTTGTAAATG ACTATAAGC TGGACGTCCA AATTGACTG CTTACCGTGA 3240  
 ACAAGCAGCT AACCTGAAA TTGAAAAAGT TGGTGAGAA TTGGGTAAAG CAATGCCATT 3300  
 CGTGTGTAAA AACGACGAT ATGCATCAA AATCTATAAC TAATTAGAAA TATATTAGCG 3360  
 TGGAGATGAT TTATGAAAA AGATTATGAG AAAAAATGCA TCGTTATTAT TGGTCTAGT 3420

208

TGTATAATG	AATTACACCG	TCGGTAATAG	TGCTAGCAGA	CCAAAATRAA	GCAGATTGGT	3480
CGTATGATGA	AAATGCTGTA	ATTAACATTT	ATGATGATGC	TAATTTTGAA	GATGGTAGGT	3540
TGCATATGAA	CTTTGAACAA	TTCTTCAAT	TGGCACAAT	AGCTAGAGAA	GAAGGTCTTG	3600
AAATTCATTC	TCCGTTTGG	AGAGCTGGTG	CGACTAAATC	TGCTCGTTAT	ATAGCGAAAT	3660
GGATTTTGAG	AAATAAAAA	CATTAACAAA	TATAGTTGGT	AAATCATTAG	GACCTAAATC	3720
AGCTGTTAGA	TTCCGGAGAAG	CTTTATCCCTA	TATTGAAGGT	CCTCTTCGCA	GAATAAATGA	3780
GACGATAGAT	GGCGTTTAT	ATCAATAGTA	GCAAAATATT	GCATCTGGAT	TGAAAGAATC	3840
GGGTTTAAAT	GACTGGACTG	CGAAAACTTT	AGCTCAGCT	ATTGCTGGGA	TATTAGATGT	3900
ACTTATTAG	GGGTGAAAT	CATATGAATA	TTACCAATTT	GTTTTCTATC	AAGACAGGAT	3960
GTGATGAJAC	TGATAGGCAA	CTGCAAAAAC	TATTTTTCAT	GTTGGATTTA	CAATTGGGAG	4020
AAATGACAGA	TCAACTAAGA	AAATTAGATT	CTAATTTTGT	TCCTCTGAGT	CAATTGTAG	4080
ACACGTTGGA	TTTGAATGAT	GTAGAATATA	AAGAAATTTT	AACTATTTT	ATCTTCCATC	4140
GTAAATGATG	TGAAGAAAGT	TTGGTAGAAT	GGTTATATGA	TTGGATTTC	ACAAATCGTT	4200
ATGAACCTCC	TAAAGAGTTT	TCGATTCGTA	TGGCTCATAA	ATACCATGAA	AGTGTACTG	4260
AACTTTTCGG	AGATGAATAA	CTAAAAAACA	GTCAATTAGTG	ACTGTTTTTT	ATAGAAAAAG	4320
AGGTTTTATA	TGTTAAGTTC	AAAAGATATA	ATCAAGGCTC	ACAAGGTCTT	GAACGGTGTG	4380
GTGTGAATA	CTCCACTGGA	TTACGATCAT	TATTTATCGG	AGAAGTATGG	TGCTAAGATT	4440
TATTTGAAAA	AAGAAANTGC	CCAGCGTGT	CGCTCCTTTA	AAATTCGTGG	TGCTATTAT	4500
GCCTATTTCC	AGCTCAGCAA	GGAAAGACGT	GAACGTGGGG	TAGTCTGCC	TTCTGCGGGA	4560
AATCATGCGC	AGGGAGTAGC	CTATACCTGT	AATGAAATGA	AAATTCCTGC	TACTATCTTT	4620
ATGCCCATTA	CTACGCCACA	ACAAAAGATT	GGTCAGGTT	GCTTTTTTGG	TGGGGATTTT	4680
GTAACTATTA	AACTAGTTGG	AGATACCTTT	GATGCCCTCAG	CCAAAGCAGC	TCAAGAAATTT	4740
ACAGTCTCTG	AAAATCGTAT	CTTTATTTGAT	CCTTTTGATG	ATGCTCATGT	TCAAGCAGGT	4800
CAAGGAACAG	TGCTTTATGA	GATTTTAGAA	GAAGCTCGAA	AAGAATCGAT	TGATTTTGAT	4860
GCTGTCTTGG	TTCTCTTTGG	TGGTGGCGGT	CTCATTTGCCG	GGGTTTCTAC	CTATATCAAG	4920
GAACAAGTC	CAGAGATTGA	GGTTATCGGA	GTAGAGCGGA	ATGGAGCGCG	TTCCATGAAA	4980
GCTGCCTTTG	AGGCTGGAGG	TCCAGTAAAA	CTCAAGGAJAA	TTGATAAAT	TGCTGATGGG	5040
ATTGCTGTGC	AAAAGGTAGG	TCAGTTGACC	TATGAAGCAA	CTCGTCAACA	TATTAAAACT	5100
TTGGTAGGTG	TGATGAGGG	ATTGATTTCT	GAAACCTTGA	TTGACCTTTA	CTCTAAGCAA	5160
GGGATAGTCG	CAGAACCTGC	TGGAGCGGCT	AGTATCGCCT	CTTTAGAGGT	TTTAGCTGAA	5220

209

TATATTAAGG	GGAAAACCAT	TTGTGTATC	ATTCTGGAC	GAAATAATGA	TATCAACCGT	5280
ATGCCAGAAA	TGGAAGAGCG	TGCCTTGATT	TATGATGGTA	TCAAACATTA	CTTTGTGGTC	5340
AATTTCOCAC	AACGTCCAGG	AGCTTTGCGT	GAGTTTGTA	ATGATATCCT	GGGCCCAAT	5400
GATGATATCA	CACGTTTTGA	GTATATCAAA	CGAGCTAGCA	AGGCAACAGG	CCCAATATTA	5460
ATTGGGATCG	CTTTAGCAGA	TAAGCATGAT	TATGCAGGTT	TGATTCGTAG	AATGGAAGGT	5520
TTTGATCCAG	CTTATATTTAA	CTTAAATGGT	AATGAAACGC	TTTATAATAT	GCTTGTCTGA	5580
GGACTAATAA	AAAAATATCA	TACCTTCATT	TTGATTTCCT	ATCTATTGAC	AAGCATATGC	5640
ACACTGTCTT	TAATACTCTT	CGAAAATCTC	TTCAAAACCA	GTTAGCTCTA	TCTGCAACCT	5700
CAAAACAGTG	TTTTGAGCAA	CTTGGGGCTA	GCTTCCTAGT	TTGCTCTTGG	ATTTTCATTG	5760
AGTATAAGT	ATGATTGAT	TTCTTTTGT	TGACAAATAT	ACTATATTA	AAAGATATAT	5820
AAGTAATTA	CTGAGCTTAT	CTGTCTTGTC	ATCTCTATTA	AGGATGGTTT	AGATAATCGG	5880
GTCTCTGCTT	CTAGGCTAGC	ACCTCAATAT	CCAAAGGAGT	GATGAATTGG	AAGGACATAA	5940
GGAATACCTA	TCTCTCAGAT	GATTTATTGA	GGAAGAAAGA	TAGGAGTTTT	TGAGCTAGTG	6000
AAGGCTTGA	TTTCTAAAGG	TTAGAACTAT	CATCTTCAGT	TCTTAAATCG	AAGAAATAAG	6060
CTATCTTACG	GAAATAGAGA	AGCATTTTTT	AAGAACTTGA	ATAATTTCGC	ACCTTAAGAG	6120
GGTAATAATA	CAGTATTTTT	ATTAGCAAT	ATTTATGGTG	TAGAGGCTAG	CAAAACCTAT	6180
ATATTATCGG	ATTTAAAAAG	GAAATAAGAA	A			6211

(2) INFORMATION FOR SEQ ID NO: 9:

- (i) SEQUENCE CHARACTERISTICS:
  - (A) LENGTH: 7939 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 9:

CGGAGTCCC	CACGATCTTT	CAAAATAACT	GAGTATATTT	CTATCTTGAT	TTTCAGATAT	60
AAATTCCTCC	TTCTGTGGCC	TCTTCTTACG	CTTGAGAAGA	GCTTCTCCGA	CATGGCTTCT	120
TCCTTACTGA	GCAAAACCTT	GAGCATAGAT	AAGTTTGACT	GGCAAGCGTG	CTCTTGATATA	180
TTTGCTCCC	TTCCCACTAT	TGTGGATAGC	GAGGCGTCTT	CTCATATCAG	TGATATAGCC	240
TATATAGTAG	GATCCATCAC	GACACTCCAG	AACGTACATA	TAAGCCTTAT	GATCCATAAT	300
AAATCTCTTC	GATTTGGGCC	GTATAAGAGC	CATCATCATT	GTGACAAATC	AAAGGAGGTA	360

210			
AGACCTTAAA GCCACPTGTT GAGCCATCCT TGATCGCCTC AATCAAAGC ATATTGGCTT	420		
CCTTTTCTCT TTTTGATATA ACAAACTGCA GCGCTTAGG GGCTAGATTA TGTCGTTTTA	480		
ACGTATCCAA AATATCCAGA AGTCGATCAG GACGATGAAC CATGGCCAAA CGCCCATTAG	540		
ACTTGAGAA ATCTCGGGCA CTACGACAGA TTCTTCCAA ATTAGTCGTG ATTTGGTGTC	600		
GAGCAAGAG ATAATGTTCA CTCTCGTTCA GATTAGAATA AGGATTCACC TTGAARTAGG	660		
GCGATTACA CAAATCATA TCCACCTTAC TCCCTTGAA GTGAGCAGGC ATATTTTCA	720		
AATCATCGCA GATGACCTGC ATTTGCTCCT CTAATCCATT CAAACGGACA GAGCGTTCAG	780		
CCATATCCGC CAAACGCTCC TGAATCTCAA CAGACAATAT CTGTGCTTGA GTACGAGTGC	840		
TAGCAAAAG CCCACTGCT CCATTCCAG CACAGAAATC CACAATCAAC CCTTCTTAG	900		
GAAAACOTGG AATCGTGAT AAGAGAACAC TATCCACGA ATAGCTAAAA ACCTCTCTAT	960		
TTTGAMGAT TTTGATATCT GTCGAAAAGA GCTGGTAAAT GCGCTCTCCT GATTTTAATA	1020		
ATTGCTTCTC TTCCATGGTC CTATTATAGC AAATTCATAT TAACATTACA AAAAATATA	1080		
AACTCTAAAC TACTTCTTCT TTTTAAATG GTGCAGGGCT TCTCCAGTCC AGATTGGTAG	1140		
CATTGTCGA AAGGGAGCAA AGCGTAGTT AAAGCGGTCG CTTGAJAJGC GTCTCCGTCT	1200		
AGGAAACTGG TACTTTTCTT CTCCAAAGT GCGGATAGAA AGACTGGCTT TCCCTGTAAA	1260		
TTCACTTAAA TCCACTACCT GAACCTGAAC CTCTTCATCG ACTTCAAGG TTTCAATGAAT	1320		
ATTTTCAATA AATCCTGTCC GAATCTCTGA AATGTGAATC AGCCCCGTAT CACCGTCTC	1380		
TAACCTAACA AAGGCACCGT AGGGCTGAAT CCTGTATAA CGCCCCTTA GCTTATCACC	1440		
GATTTCATC TTAGTCTCG ATTTCAATAG TTTCATTAAC AACATCTCA ACTGGCTTGT	1500		
CCATAGCTCC TGTCTCAACA GCAGCAATGG CATCCAAGC AGCGTAAGAT GCTTCATCAG	1560		
CTAATCTACC AAAAACGGTG TGACGGCGGT CTAGGTGAGG TGTCCCACCT TGATTGGCAT	1620		
AGATTTCTGC AATCGGTCTT GGCCAACCAC CACGATTAAT TTCTTTCTTA GAAATAGGTA	1680		
GGTGTGGTT TTGCACGATA AAGAACTGGC TGCCCTGGT ATTTGGACCA GCATTTGCCA	1740		
TGGAAGAGC ACCACGGATA TTGTAAAGCT CTCTGAGAA TTCACTCTCA AAGATTGCC	1800		
CGTAGATTGA CTCGCCACCC ATACCAGTTC CAGTTGGGTC TCACCTTGG ATCAATAAGT	1860		
ACCTGATAAT ACGGTGGAAA ATGACACCAT CATAGTAGCC ATCTTTTGAA AGAGATACAA	1920		
AGTATGCCAC TGTTTTAGGA GCATGTTTCA GGAAGAGCTT GATACGTAAG TCTCCGTGAT	1980		
TGTCCTTAAT AGTCGCAAGA GGACCTTCTA CTGTTTCAAT GTCTACTTGT GGAATAATGA	2040		
ATTCCTTTTC TACCATACCA AATACTTCTA AGGCAGCAAA AATGCCATCT TCTTCTAATG	2100		
TTTTTGTAAT ATAACTTGCT TTTTCTTGA TTTTATCATG AGAAATTCCT ATGGCAACCC	2160		



TGATTCCAGC	ATAATCAAAG	AGTTCCAAGT	CGTTGAGACC	ATCTCCAAAA	ACCATGACCT	2220
TCTCTGGTTT	CAGGCCAAGG	TGTTCCACAA	CCTTTTCAC	CCCCGTCGCT	TTGGAGCCCTG	2280
AAATCGGCAC	AAATATCAGAC	GAATGTTGAT	GCCAAACGAAC	CAATGCGAAGT	TTGTCTGAGA	2340
GACTGTCAGG	CAAGTCCAAG	TCAATCTCCCT	TATCTTCAAA	AGTCCACATC	TGATAGATAT	2400
CTTCTTTTTC	ATGGAATATCG	GGATCTACAT	CTAAGTCGGG	ATAAATTGGA	TTGTATAGCTT	2460
CATCTATCAT	ATCGGTGCGA	GTGACAACT	TGGCATCATG	ACTCCCAACC	AAGCCATACT	2520
CAATTCCTTC	TTGCTTAGCC	CAAGAGATAT	ACTCCTCAAC	ATCTGACTTT	TCAATCTGAT	2580
CGTGATAAAT	GACCTGACCT	TTTTTATCTT	CGATATAAGC	CCCATTCAAA	GTTACAAAAA	2640
AGTCAGGCTT	GAGATCACGA	ATCTCTGGAA	CAACACCAAA	AATGCCACGT	CCAGAGCGGA	2700
TTCTCTGTAA	AAATTCCTTT	TCAACGCACT	GTTTAAAAAC	AGTGGGAATT	GTAGTTGGAA	2760
TAAACCTGT	CTTGAAATC	CGCAATGTAT	CATCAATATC	AAAAAGACA	ATCTTGATCT	2820
TCTTTGCCCT	GTATCTTAAT	TTGCGCTCCA	TCTCACTACC	TCTTTCAATC	TAATCTTTTC	2880
CATTATATCA	TAAAGTAGCC	AAATCCCCCTA	TTTTTCAAAA	GTTTATCATT	TTTATTTTAA	2940
TTTCTTGGAT	GAGAAAGAG	ACATATTTAT	GAAAAAGCTC	CATCGTGCTT	TAAATGTGTT	3000
CTCTGTGTTT	CAAACTCGTA	AAAAGGGAGC	CAGTATCTCT	AATCTGCTCT	CTCATTTCAA	3060
AGCTTTGTGAA	AAAAGACCCG	TTGGGGTCTT	AAATTCGCTT	CTTGTTTTCA	AGCTCATGAA	3120
AAAGAGACCC	AATCGGTCT	TTTCTTTAAT	CTTCGTTTAC	GAAAGGCATC	AAAGCCATTA	3180
CGCGAGCCG	TTTGATAGCT	GTGTGTACTT	TACGTTGGTT	TTTAGCTGAA	GTTCCTGTTA	3240
CACGACGAGG	AAGGATTTTC	CCACGTTCTG	AAACGAAAGC	GCTAAGAAGC	TCAGTATCTT	3300
TGTAATCAAC	ATATTCAAAT	TTGTTTTGCTG	CGATGTAACT	AACTTTTTAA	CGGCGTTTGA	3360
ATCCGCCACG	ACGTTGTTGA	GCCATGTTTT	TTCTCCTTTA	TAAAGTTAGT	TGTCATATAG	3420
AATGGTAAAT	CATCATCTGA	AAATTCCAAT	GGGTTTGTGG	CTCCAAATGG	ATTTTCATTA	3480
CGTGAAAAAGT	CTGGTACTGA	ATTTGTAGGT	GCTGAATAGT	GTGCACTTGG	TGCAGAGTAA	3540
GCTCCACCTG	TGTGACCCCTC	ACGCACACTA	CGGCTTTTCCA	ACATTTGGAA	ATTCTCAGCC	3600
ACGACCTCTG	TCACGTAGAC	ACGTTGTCCT	TGCTGGTTAT	CGTAACTACG	AGTCTGGATA	3660
CGACCTGTCA	CCCCGATAAG	TGAGCCTTTT	TTAGCCCACT	TAGCAAGATT	TTCAAGCCTGT	3720
TGGCGCCACA	TAACGACATT	GATAAAATCA	GCCTCACGTT	CACCATTTTG	ACTCTTAAAT	3780
GTACGGTTTA	CTGCAAGAGT	AAAAGTCGCA	ACTGCTACAT	TTGATGGGGT	ATAACGCAAC	3840
TCAGCGTCAC	GTGTCTATAC	CCCTACAAGT	ACAACATTGT	TAATCATAGT	TTACCTTCTT	3900

		212	
ACGCGTCAAT	TTTGACGATC	ATGTGACGAA	GAATGTCAGC
CAAACTCTTT	AAGAGCTGCA	TCGTCAATTG	CTTCAACGTT
CACGGAAATC	TTGGATTTCG	TATGCAAGAC	GACGTTTTTC
CAGTTGCACC	GTTGTCAGTC	AAAATAGAGT	CAAAACGTGC
CTTCTTCAAT	GTTTGGACGA	ATGATATAAA	CAATTCGTA
TCCTTTTGGT	CTAATGACCC	CAAGACCTTG	CAAGGGGTAA
ACTATTATAC	TAGAAAAAAT	TTTTTTACGC	AAGTAAAAAC
CACATGGGCG	TTTTCTCGTT	CTTATGGTTT	GATACGGTGC
AGCTTCAAGG	ATATGTTTTC	TTCTGCTGTC	GAAGGTTACC
AAATCCTCCG	TGTGGAACTG	TACCGTATTT	ACGAAGGTCA
ACGATCCATG	CCAAGTTGAT	CCATCTTAGC	GACAAGGGCA
AGACCCACCG	ATAATTTCTC	CATAGCCTTC	TGGAGCAAGC
CTCTGGATTT	CCAGGAACCTG	GTTTCATGTA	GAAGGCCTTG
GACAAATGTT	GGCACACCAA	AGTGGTTTGA	AATCCAAAGT
ATCACCATGC	TCAAGATGCT	CGTAGTCAGC	ATCTTCATCA
GTCAMTGGCT	TGATCGTAAG	TGATACGTTT	GAATGGCTCT
TTCTGTATCA	CGTTCCAAGG	TTTCCAAGGC	TTGAGGCGCG
AAGAGCTTTC	ACATAAGCTT	CTTGCAAGTC	AAGCGACTCA
CTCAGCATCC	ATCATCCAGA	ACTCAGTCAA	GTGACGGCGT
GAAAACTGGA	CCAAAGTCAA	AGACACGACC	AAGAGCCATA
CTGACCTGAT	TGGCTCAAGT	AGGCTGGCGT	TCCGAAAGTAG
AGAATCTTCT	GCCGCATTTC	CTGAAAGAAAT	TGGGCTGTCA
GTCAAAGAAC	TCATAAGTTG	CATAGATAAT	AGCGTTACGG
CTTAGAGAGG	CGTAGCCACA	AGTGACGGTT	ATCCATCAAA
TGGTGTGATT	GGTAGTCTTT	GAGATTCCAC	GATCACTTCG
ATAGCCCAAT	TTAGAACGTT	CGTCCCTCTT	GACAATACCT
TTGGCTCAAG	CGTTTGATAA	CATCAAACTT	CTCAAGTCCC
GACAAAGTTT	GGTTTAAAGG	CCACACCTTG	AAAGAAGGCT
GAAAGCGATT	TTTCTTTTTC	CTGATTTGTT	GGCAACCCAA
ACCAACATAG	TCTTTTACGT	CAATAATCGT	TACACGTTTT

3960  
4020  
4080  
4140  
4200  
4260  
4320  
4380  
4440  
4500  
4560  
4620  
4680  
4740  
4800  
4860  
4920  
4980  
5040  
5100  
5160  
5220  
5280  
5340  
5400  
5460  
5520  
5580  
5640  
5700

TTTTTATCTT TTATGSCAAA CCACCTCTAT ATTGTTCCTCA TCCAGGTCAA TCATAAAAGC	5760
AGCATAGTAA ATCGGATGCT CACTTCGATA ACCAGGAGCC CCATTGTCTC GCCCACCTGC	5820
CTCTAAGCCA GCCTCATAAC AAGCCTGAAC TTCTTCCTTA TTTTCTGCTA AAAAAGCAAA	5880
ATGAACAGGA TCTTGTGTTT CTTGAGTCAG CCAAAAATCA CCACCAAGAT GAGGGCTGTT	5940
CGGGATAGA AAACATAATTA GAGAACTAGT CTTAAAAGCC AATTATAGT CCAAAGGAGC	6000
GAGAAACTC CTATAAAATC CTTATGAAT TTGTAAATCC TTTACCTTAA TCTCAAAATG	6060
ATCAATCATT CTCACTACCC ATAAATGCTT TCAAGCGTTC GACTGCTTCT TTAAGCGTGT	6120
CTAGGCTGTG UGCATAGCTG AGGCGGACAT TTTCTGGTGC TCCAAATCCA GCTCCTGTTA	6180
CCAAGGCCAC TCGGCTTCT TCTAAGATAA CAGTTGTAAA GTCTGTACA TCCGTGTAGC	6240
CTTTCATCTC CATGGCCTTT TTGACATTTG GGAAGAGATA GAAGGCCCTT TCGGTTTGA	6300
CCACTTCAA TCCTGGTACC TCTGCAAGGA GGGATAGAT GGTATTAAGA GGTCTCTCAA	6360
AGGCTGTAGC CATGCTTTCT ACAGTATCTT GTCACCTGA TAGAGCCTCA ACTGCTGCAT	6420
ATTGGCTAC TGCTGACGGA TTCGAAAGTTG TTTGACCTGC AATCTTGGAC ATGGCAGCGA	6480
TAATGTCTGC TTCTCCAAG GCATAACCAA TCCGCCAACC AGTCATGGCA TAAGTTTTAG	6540
ACACACCATT GATGACCACT GTTTGCTTGC GAATCGCTTC CGATAGGCTA GAAATCGGTG	6600
TGAAGCTAGC ACCATTATAA ACCAAGCGGC CATAGATATC GTCTGCTAGG ATGAGAATAT	6660
CATTTTCTAC AGCCAGTTT CCAATTGCCA AGAGTTCTCT ACGGGTGTA ATCATACCTG	6720
TGGGATTAGA TGGCGAATTC AGCACCAAAA CCTTGGTCTT GTCAGTGCGA GCTGCTCTTA	6780
ACTGCTCTAC GTTCACCTTA AAGTGATTGT CTTCCTTAGC AGAAACAAAG ACGGGAACGC	6840
CTTCTGCCAT CTTGACCTGA TCTCCATAGC TAACCCAGTA TGGGGTTGGG ATGATGACTT	6900
CATCACTGG ATTGACCAA GCCATAAAGA AGGTATAGAG AGAATATTG GCTCCCGCAG	6960
CGACTGTAC TTGATTGAC GCTACAGAAT AGCCGTAAAA GCGCTCAAAG TAGCTATTGA	7020
CCGCCGCTT AAGCTCTGGC AGACCTGAGG TTACTGTATA AAAAGAAGCA CGCCATCTC	7080
GAATCGATGC AATGGCGGCA TCTTGATAT TTTTGGAGT AGTGAATCT GGCTCACCCA	7140
AGGTTAGAGA CAAATATCT CTACCCTCAG CCTTCAGTGC TTTGGCAGG GCTCCAGCAG	7200
CCAAAGTCAC ACTTCTTCC ATTTCTAAAA CACGGTTGGA TAGTTTCATA GGCCCTCCTT	7260
GTTGACCAAT GCTCCTGTTT CAAAATCTAC TAGATAAAAA TCAGATCCTG ACTTAACCTC	7320
CCAGATTGGC TTATCTTGAT AACGGCCAAA GGTATCTTTG TCAATCTCGC CAGCTCCCTT	7380
TTCTTAGTAA ACGGTTCTG CTTTTCCTTG TGAACACCC TGATTTAGCT GATTAACGTA	7440

214

AATCTTATGG TCATCTTAC CAATCAGGAC AGCAAGCGCT TCTTGCCTGT TGTTCAGACC	7500
AAGAAGCGCTG TAATAAGATT CCAAGCCATT GTATAAATCA ACCTGATCAG CCTGCTCTAA	7560
TCTTGCATAC TGCTGAGCTA ATTTTCTCC TTCACTTTTA GCTGTTTGAT AGGGTTTCAT	7620
GCTAAGAGAA ACCATATACA GAAAGGAACU ACTGATAACC ACAAACAAA TCGTCATCCC	7680
TAGACCATAT TGCCACAGTA GATTATTTTT TGCTTGTGTT TGTCTTTTTT TCACTGCTCT	7740
ATTTTACCAT CTATTAAGCT TTATTACAAG TGAATATAAG AATACTCTTC GAAATCTCT	7800
TCAAACCACG TCAGCTTTAT CTGCAGACCT CAAAGCTGTG CTTTGAGCAA CCAATTCAT	7860
TTCTCCCTTC AAACAAAACC GATTTTGAAA GTGAACACGT TCTTACTTTT TCAGTCACAA	7920
ATGATTAGAG TTTGCGXGG	7939

## (2) INFORMATION FOR SEQ ID NO: 10:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 9897 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 10:

CCGCTCTACC GTCAATAAT TACCATTGTT TTTAATACCG AAATTTTAT CTACTGAAAA	60
TTTCACTGGT CTGTTGGTAC GATCGTCGTA TACAGTACCA TTCTACGAA TAGTATAATT	120
GTAATCAGTA TCACCTTGTT TCCTTAATTT AAGTAATAAA TTACCATCAA TTTGTTTATA	180
ACCTGAATCT TTTCTAGTTG CTTCCTTAAA ACTTACTCCA GCAGGCATCA CATCAGCAAA	240
CATGAGTACT TTTTGTGTTT TTTTTCACAC AATAACAGAG TCAATATAGG TTGCAACCAC	300
GCTGATTTGT AAGTCACGTC CACCAACTTC ACGAGGCCMT TCTAATGCTA CTGGCGCAAA	360
ATCATCGAAT GCCAATGTTA ATTTTGGTTT AGTCCATGTC TTACCATPAT CATCACTATA	420
ACTTGTAACA ATATTAAATTT TATTCAAGAA ATCATGAGTT CCACCGTAAC GAGCGTCAAT	480
GCTTGAAAT ACCCGCCAT TGCTAAAAGT ATACGAACAT GGAATACGGA AATAGTTAGA	540
ACCTGTGCTA TCATTAGCCG TATAAATPAA ATGTCCAGTA ACAGCGTTTG TTGTCATCTT	600
TTTAACTGTT TCTTCTATCCA ATGCACATTT AAAGAATTTT ATATTTCTTA GTGTTCCGTT	660
AAAACCAAC GCCGTTTTTC CTGCACGTTT CACTCCGCCA AGCATATAGT AATCAATACC	720
TTTAAATATC TTGATGTTTA GGAATATATC CACTTTCTTT TCTACTACTT TTGTACCATT	780
TGCGTATPAA GAATATGTTT TTTTGACTGA ATCTGCTACT ACTGCAACAG TGTGATGCAC	840
AGGCTCTTGT TTGTACTTAC CCCAAACTGA ACGAGGCTGT GATACTAGGT TATTTTATTT	900

GGAGAAGTA TCACGCGCTT CCAATCCCAG CTCACCATTTG TCTCTAAGGA ACACATCTAC 960  
ATAACTATT TGTTCACCGG GTTTTGGAAAT AGATATTCCA AACAGAGCTT GTAAGCCCTT 1020  
CTCACTTGAC TGATTGTAAT TAATCACTAC AGTAAAGTCA CCGCTAGTAA ATTATCCCTT 1080  
TAACCTCTTA GTAACATTTT CTCGCCCCC TGTAAAGTA ACATTATTTT TTTCTAAGAC 1140  
AGGAGTTTCT TCCGCTGTAG AAGATGGATC CTTAACAGTA GTTCAACTG TFCGAGGTTG 1200  
TACAGTAAT TCCGAAGAGT TATCCGATGT AGGTTGTAAT TCCGAAATCG GACTCCTTGG 1260  
TGCAACAGGT TGCACCAAT TFGGTGTTGA TACTTCAGAA GTTTCAGTCT CCGTAGCTGC 1320  
AACTGAGTTA GCAACAAATG CTGATAATAC CACTACAGTA CCTAAGTTA CATATTGTTT 1380  
AATATTTTTT TTCAATTTAT TTTTCCTCGT TTAATACTTT GATAACAAGT TTTTAAACAG 1440  
TTTCATCAT TGCATGAATC TTTGGTTGGT GAAGATCTTC TTCAAAAGTC ACCAACATAT 1500  
TCCCCTGAAG CAATTCACAA ATTTGATAGT CTTCGCTATC GTAAAAAGCA ATATCCTCTT 1560  
CTTCGCTTAA AGGTACACCT GACTGGGCAC GAACTGGGGA AGTTACTGCC ATTTTTCAG 1620  
TATTTTCAAC AACAATATGA ATATCTAAAT ATTTCTTATG AGTTTCAAAA ATATCTCCTG 1680  
GAATCCATC AGCTAGATAA GTCAATCAAT TFGCAAAAGC ATTTTCCCGC TCAATATCAA 1740  
TTTTCCTTCA AACTAAATCT GTCAAAATTTG TATTTTCTTA AAAATCACAG ACTTTTGAAA 1800  
AATATTTATT GACAGAAGCA TATCGTTTAA AATCAGATTG TTCAGAAATA ATCATATTAT 1860  
TTTCTCTTTT CTATTAGTGA CGAACTTCCC AACTTGAATC CGCTTTAATT TCTGTAATAT 1920  
CATGAATGCT TGTATATTTA GGTGCAGATA CTTTATTTCC AGTAAGAACA GATACAATAT 1980  
AACTTGAAAC TACTGATACA GAGATTGAAA TCAATGAATA TGCCAGTAG CTAACAGCTG 2040  
TTGGAGGAAG GAAGTATTTA ATAAATACCA TGACGATGTT TGATACAATC AGCCCTGCAT 2100  
AAGCACTTTG TTTATTTGCT TTTTATGAAA CAATCCAGG AATAAATACA CCACCAAGTA 2160  
GACCAAGTAC AAGTCCCATG AAATATTTGA ACCATTCTGA TGCAGATTTA ATATCTGAGT 2220  
GAGCCATGAC AATGGAACCA CCAATTGAGA ATAAACCTAC TGCTAGAGAT ACGAATTTGTG 2280  
CAAATTTGCT ACGACGATTT TCTGACATAT TTTTAAAGAT GACATCTTGA ATATCCAATG 2340  
TCCATGAAGT TGCAACAGAG TTCAAACTG TFGAAATAGT TGATTGAGAT GCTGCATAAA 2400  
TCGCTGCCAA GATCAAACTT GTGATACCTA CTGTTAACTG GTATGCAATA AAGTACATAA 2460  
AGATTGCTC TFGAGGATA TFGTAGCTG CACTATCTGC ATTTTGTACT TGATPAGAATA 2520  
CGTACAAGCC TGTACCAATC AAGTAAAGGA CTGTTGCAAT TGCAAGTGAC AAAACACCGT 2580  
TTGTGAACAA CATCTTATTA AGTTTCTTAA TATTTTGTGT TGTAGTAAAA CGTTGAACCA 2640

216

AATCTTGAGA TGAAGCATAG GAAGACAAGA TTGTAAGGCC TGAACCCATC ACAATTAAAA	2700
AGATGGAGTT TGAAGCAAG TTAGGATCGA AAGGTTTTTC ATTTCGACGA AGGAATPTCC	2760
CGTTTGCTAA TGTTTCTGCT ACTGCACCAA AGCCACCTTT AATATTAGCA ATCAGTACAA	2820
ATAAAGCTAA AACGACCA CTAATCAGAA TCACACCTTG AATAAAGTCT GTCCATAATA	2880
CGGATTTTAG ACCACAGTA TAAGAATAAA CAATTGCAAC TACACCCATC AAAATAATCA	2940
AAATATTGAT GTCAATTCTT GTCAATCTGT ATAAACCAGC TGATGGGAGG TACATAATGA	3000
TAGACATACG TCCCAATTGA TAAATAATAA ACAAGAGTGC TGAATAATA CGAAGTGCTT	3060
TAGAAATTAA ACGTTTATCC AAGTAATCAT ATGCCGTATC GATGTCTATC CGTGCAAGA	3120
TAGGTAAGAT AAAACGAATT GTCAAGTGAA TAGCTACTAC CATCCCTAAT TGAGCAAAAC	3180
ATAAAAATCCA GCTACCTGCA TAAGAGCTAC CAGCGAGTCC CAAGAAGGAA ATCGGACTGA	3240
GCATTTGTGC AAAAATGGAT ACCGAAGTAA CATACCAAGG AACCGAACCA TCTCCTTTAA	3300
AGAATCTTTT TCTTTTCATC TCTTTTCTAG AGAATAGAT ACCTGCAACC AACACCGCAA	3360
GTAAATAAAT AATCAAGATA ATTAAGTCAA TTATTGTAAA TCTGTGTGTG CCAATAACAT	3420
ATCTCCATAT TGATTTTATT TATTATAAAA ATTCCTTTTG TGCTTGTGTA ATAAAGTCTG	3480
CTGCTTGTPT TGCAACTTCC AAGTCACCTT CTGCCAATGC TTCTAAAGGT TGACGAACAG	3540
AACCTAAATC AAGTTTTTCA TTTAGACGCA AAATTCCTTT TGCTACAGCA TACATAATTG	3600
CCTTACCTGA TATCATCTTA TAGATAACTT CATTGATAGC ATATTGAAGT TTTTATAGCT	3660
TATCTAAATC TCGTCTTGTA ATCAAACTTT CCAATTTCAA GAACAAATCT GGCATAACGC	3720
CATAAGTACC ACCAATACCA GCTTCTGCTC CCATCAAGCG ACCACCAAGA TATTGTTTCT	3780
CTGGACCAAT GAAATCAATG TAACTTCTCT CACCTGCAGC TACAACATT TGAATATCTT	3840
GTACAGGCAT AGAAGAATT TTAATCTCAA TCACACGAGG ATTTTGACGC ATGTGTGCAT	3900
ACAAATACCC AGTCAACGCA ACCCTGCCA ATCTGTGAAT ATATTAGATA ATAAATCTG	3960
TATTTGACGC AGCTTCACTC ATTGCATTCC AATATGCTGC GATTGAATAC TCTGGCAATT	4020
TGAAATAAAT AGGTGGGATA GCTGCAATAG CATCGACTCC AACACTTTCT GAATGTTTTC	4080
CCAATTCGAT ACTATCTTTC GTGTTATTAC ATGCAATATG GTTGATAACT GTTAATTTAC	4140
CTTTAGCAAC TTCCAATAAC GCTTCAATAA TTTGTTTACG ATCTCTTACA CTTTGGTAAA	4200
TACATCTACC TGAAGAACCA TTTACATAGA TACCTTTTAC ACCTTTGTCA ATGAATATAT	4260
GTACACGAGA TTTTACACGA TCTTGGCTAA TTTCAACATT TTATCATAG CAAGCATAAA	4320
ATGCAGGGAT AACGCTTTTG TATTAGTTA AATCTTTCAT CAGATTCTCT CTTTATATTG	4380
TTTTTATTAT GATGACATTA ATAAATCGCT GAGCAATTTT TTTTGGAGCT GTAAATCGCTC	4440

CACCAATGAC TACACTGEEA ACACCTAAAC TATAAGCTTT TTTTAAATGT TCTGGATAAT	4500
GAATTTTCTC TCGGCAATTA CCGGAATATT AAAATCAGCC AATTTTTTCA TTAGTTCAAA	4560
ATCAGGCTCA TCTGATTGTA CACTTGTACT TGTGTAACCT GATTAATGTTG TACCAACAAA	4620
ATCAACGCTT GATTTAAATG CATAGAGACC TTCATCTAAA TTACTTACAT CCGCCATCAG	4680
CAATTGATTC GGATATTTTT CTATTATTTT TTTGATAAAT TCACTGACAA CTAAGCCATC	4740
ATATCTTGCT CTAAAGTTG CATCAAATGC AATGACTGTT GTTCCGCAAT CTACAATTC	4800
ATCTACTTCT TTCATCGTAG CAGTAATATA TGGTCTTGA GGTGGATAAT CCCTTTTGAT	4860
AATCCCAATT ATTGGTAAAT CTACTACTTT CTGAATGCTT TTAATATCAC GCACAGAAAT	4920
TGCGCGAATG CCCACTGCTC CTGCTCTTAA AGCTGCTTTA GCCATAAAG GCATCAAGCT	4980
AAATTCTTCA TTATAAAGGG CTTCAACAGG TAAAGCTTGA CAAGAAACAA TGACTCCACC	5040
TTGAACCTGG CTTATAAATT TTTCTTTAGT CCAATTTGG CTCATTTTAT TATTCCTCCT	5100
TATGGATAAT AGTTTGATTG TAAATAATAT GTCTCTCTGG ACTTCCAGA TAATTAGAGA	5160
ATAAGCAGTC TGTAATTAAA AGTATTGGAA ACTGAGGTGA TATGCGATTG CCATACGAGA	5220
GATGATCGGT CGAAGCTAAT AACAAATAGT CATCAAGAA ACNATCTTCT TCGTCAAAT	5280
TTCTTGTAAG CATTAAGACT GTTTTAGGCG CTTTATCTGC AGCTTTTGT AGACCTTCTA	5340
GTACAAATTC AGTTTGACCT GAATGGATG CTCCAATGAC AAGGCAATTT TCATTAAGTA	5400
GTAAGCTACT CCACAAAATC ATATCCTGCT CTGATAATAC TTCACCAATC ACTCOAGAC	5460
GCATAAATCT CATCTTCATT TCTGTAAAG CAAGAACAGA ACTTCCTTTA CCGTAGAGAT	5520
ATACAGCTTC AGCAGTTTCT ATCATCTCAG CAATAOGCTC AAGTTGAATC TCATCAAGAA	5580
CGGTGAAAT TTTTCTCAAC ATTTCTCAT AGTCGGATAA AACTTTTCTT GTTGCTCTG	5640
TATATAATGC CAACTTTTCT TTCTCATGAA TCATCTCTTG GTATTGAAA ATGAATTGTC	5700
TAAAACCTTT AAAACCATAT TTTTTCGCAA ATCGAGTCAA TGTGCTTTG GATACATTAA	5760
GGTATTCGCA CAATGCTTTA GATGAATAAT CATTCAGAGG TTGCTGTTT AAGAAGAAAT	5820
TAGCAATGTC TTTTTCAGCA TATGCCATAT TTGGTAAGTT AGCTTCTATC ATTGGAATTA	5880
GTCTTTTGTG CAGTAACATA TGAGCTCCTT AGTTGAAGTA AACGTTTACA TTCTTTATTT	5940
TAAACACTTT TTTTTTTTC AATATTTTTC ATAAATTAGA AACTAGTTTC CAATTTCTTT	6000
CGTTTCATAA CAGAACAACA AACATAAAAA TATAATAGTT TTTATCTTTT TTATCGTAAT	6060
TATATGTATT GTAAGAACGT TTATCACTAA TAATATGPTC ATATPAAAAA ATTTTAGTAA	6120
TATTTTATTT TGGTTTATTT ATTTCTTTTC GGAATTTCTA TATAATATTT TATTTCTAAA	6180

218

AAAATTGAAA	AAATATTCT	AGTTTCCTTA	TTTTATATAG	GTAATATATT	TTATTTCATA	6240
ATTAAAAGAG	AATCCCATAA	AAACTACAGA	TTTATGAGAT	AAATCAGGTC	ACCTATTPTTA	6300
AAAAGCAGC	AAACTATAAA	CTAAAAGTT	CCACACCAAA	TGFAACCCCA	TACTTCCCCA	6360
TAAGTCAGAT	TTATAGCGCA	CCATACCTAA	AARCAATCCA	AGTGAACGT	ACAGACACCA	6420
AGCTAGAATG	GTTCTGGAT	GATGTACTAA	GGCAAAATAA	ACACTTGTCA	AAGCAACTCG	6480
AATACTAAT	TTTCTAACCA	AGTTCCATAA	AATTTCACGA	TACAGAAATT	CTTCAACCAT	6540
ACTGCAATG	ATTAAAGAACA	ATAAAATGA	AAACCAAGGA	ACTTGATGTT	GAAGGCCAAT	6600
TAAATTGTGT	TGATTGCTGC	TTCCCTTGAGC	ATGAATCAGG	CTAAAACATA	GACTTATAAT	6660
CAGTAGACTA	GCTAGTCCAA	TACCAAGGCA	TTTCATCCTA	GTTTTCATAT	TGACCTTGAC	6720
CACCTGTTTT	CGTTGACCAT	ACATCCATAA	AAAAGAAAAA	AGAGACGCAC	CATAGAGAAC	6780
CTGTAGTATA	GTTAACTCAC	CGATACAAGG	AAATTTCAAT	AAGTATAGAG	ATACCAATAG	6840
GACATTACTT	TGTTGGAATA	TATAAACTGG	AATTATTCTT	TTCATAGTTA	CCTCCGAAAT	6900
AAATCTTCAT	AATCTAAATC	TAATATCTGC	ACAATCCTTT	CTACCATGG	ACTTTGAGGC	6960
ATTGCTTGTT	CCATCTTGTA	GTGGCGAATC	TTTTGATATA	AACGATTCAA	TTCACCTTGA	7020
TAGTGAAACT	CTCCCGCAAA	CATTTTCTCG	GTTAACTCAA	TCCAGCTGAT	ATTCTTTTCA	7080
GCCAAAATAA	TGGACAAGTT	CTCCCAAAAT	CGTTCAAGCA	TATTCTCTCT	CCTTTAGTTA	7140
GATAAATAAT	GTGTTTGYGC	CATGTAAATC	AATTGTTTCG	TATCTCTTGG	CAATAGAGCT	7200
CTAGCCTCTT	CCAAATTCAG	ACTTGGATAA	ACCGCTTAT	TTGAAACCAAC	AAAAGGAAGT	7260
CCGATGGTTA	GTTCAAGGAT	TTTTAAAAAT	ATCTCAACGA	AATCCGTTAA	TCTTAGATTG	7320
TCACGGTTCT	TAAATCGTAA	TAAATTGGGA	GATAAAAACT	CAAAAACATC	TGAAGAATAG	7380
CTCATCATCT	CAATTAATTT	GTCCCTTGTC	ATTTCAAGAA	CTGAATGACA	AGATACCTCA	7440
ATGCCATAGT	TTTGGAGAA	GTCTAAAAGA	AGTTGATTTT	TTTGGCTATT	TTTACTTAGA	7500
TAGAGATCAA	TCATGGGAGA	CCTCCAACAA	ATTTGCTTCC	ATTTGATATT	CTGAGACGAT	7560
TAAGGAATCT	AACAACCTTG	AGAAGTTAAT	CGATTTCTTG	TCTTCATCAT	AAGCTTTTAC	7620
AGTTACTTGG	GTTTAAAGTA	TCCCTCTCTT	TCCCTCGGCT	CGATAGTCTT	GTCANATATA	7680
AACAAAAACA	AGATTCTGAT	TATCATCTAC	AAAGGCATTA	ACTCCGTCTT	TTATATCCTG	7740
ACTTCAAGG	AATTCCATAA	CGTTTGAGAG	ATAGGATTCA	TAAAAATAGT	GGTAATTATG	7800
TTTTTTATGG	TAATCATCTA	AAAATGTTAC	CTCAAACTCA	CATGGATAAT	TGGGCATCAA	7860
AAATATTGTT	TCATCCAGCT	GTTTGATTTT	TGCATCATGT	AATTCTGTTT	CTAATTCATC	7920
ACAATCTAGT	ATTGATTCTT	TATTTAATGC	TTTTATCTTT	TTCTCTTATT	TCTTTTAATT	7980



TCCTTGGGAT	TGCGGCATC	ACAGGAACGG	TTACACFAT	ACCAACTTGT	TTATAGAGCT	8040
GACTATTAA	AGAGACTTTT	CTAGCAGCTT	CAAAAGCCTA	ATCAGGAAAG	CCATGCAATC	8100
GAAACACTC	TTTAGGAGTG	ATTGCTGCTA	TTCTCAAAAG	GTAATAATGT	CCATCTATTA	8160
AAACACCAGC	TACTTGGTAA	ACTTGTATAT	CTTCTCCTTC	ATAGCTAGCC	ACTACTACTC	8220
CCATTGAGCC	ACTAGTGTGT	AACGTATTAG	CTATACCTTT	TCCAACCTTA	CCACGACGAT	8280
ACTGAGAACT	TGGCTCTTCT	AATTTGATTT	AATCCCAAT	CTCTGCTTGA	GCATATCTCT	8340
TTTTGCTTGC	TTCCCGTACT	TTTAGAAAT	GGATTGGTTC	TGGAATTAGT	ATTTTGGGGA	8400
TTTATCTTCC	TCTTGCATC	GTAGTCAGTG	TTGGAGAATA	GCCCTCCTTT	CCATAGACAC	8460
GACCTCTCTC	CTTAAAGCTA	GTCGGTAAAT	CTCCAACAA	GACAATGCCA	TAACGATCTCT	8520
GAGTATTFAA	AGTAAACATC	GGCTCTTGAT	TTTCTCTTAA	GCGTCTCCCA	TTTTGTCTCT	8580
TGTCTAATCT	ATCTGTGCTC	ATACAAGGAA	TGCGAATCTT	AAATCCTTCT	CCCTTACCAC	8640
GAACTAAGGT	TGGCGCAAGA	CCTTCTGAAT	AATAGACTTT	ACCGCTCAT	CCACTCTTTG	8700
ATGGATTCAA	ATTCTCTAGT	GCTTCAAAAG	CTTCAGAGTT	AGTTGCTTGA	CCCTCTCGTC	8760
TGAAAGGAAA	TAAGAGGCTG	GTAACCTTCT	TTCTAGAGAT	TCCGATAATA	AACACCCCTC	8820
CTCTGTTTTT	GGGAACGCCA	AAATCCTTAC	TGTTAAGCAC	CTGCCACTCA	ACATCAAAAC	8880
CCAACTCATC	AAGTCTGGTA	AGTATTGGGG	TGAACCTCCG	TCCCTTATCG	TGATTGAGTA	8940
GGCCTTTAAC	ATTCTCAAGA	AAAAGAAAAC	GTGGTTGGAT	TTGTTTGGCC	GCCCGAGCAA	9000
TTTCAAGGAA	CAAAGTTCTT	CTAGTATCTT	CAAAATCCAA	TGCTTCTCCT	GCGATTGAAA	9060
ATGCTTGACA	AGGGAATCCC	CCACAGATGA	CATCGACTTT	CCCTCTAAGT	TTTTTAAATT	9120
CGTCACTTGA	AACATCTCTG	ATGTCATGAA	ATTCTATTTC	TCTTCTCCGT	TGAATAATGG	9180
ACTTATAAGA	TTTCTTAGCA	AATTTATCAA	TCTCAAAAA	TCCCAAGCAC	TCAATGCCCT	9240
GAGCTTCCAT	TCCCATCTTA	AAGCCTCCTA	TCCGAGCAAA	TAAATCTAAA	ACCCAAATCA	9300
TTTATACCTC	TCTCACTTAG	ATGTAACCTA	CAAAACCCCT	GACCTCATGA	GCCACTTTCT	9360
TCTCTCTCAT	GAGGTGAGTT	TTACTTTCTG	CTGTCTCAGT	ATCGTTTCTT	CTCGATAGAT	9420
TTCTCTAAAA	GGGCAGACTC	CTCCCTTGGT	TGCTCACACG	ATTTTCTTCA	CTCGACTGTG	9480
CTTTAATGCA	TCAATTACGA	CGCTTTTCTT	CTAGGTGGTT	CATTAAGGAAC	AGGAAGATTTC	9540
AGGTTGACTT	TTCTAATCTT	AGAAATAAGT	GCTGAAJACA	ATTTCGGAATA	GGCATAGAGA	9600
CTAGACAATT	TGAGGAGCTG	CTTGGCTCTT	GTTCGAACAC	ATTTTCTCTAC	CACGTGAAGA	9660
AAAGAGTGGC	GGAAAGCTTT	GATTGTATAA	GTTTGGAAGT	CACCTCCAGC	TAGATGTTTTG	9720

220

AGAAAAAGAT AGAGATTGTA GCGATACAG CTCATCATCA TACGAACCTG TTTTGTGATTA	9780
AGGTTGAACCT ATCCGTTTGA TCGCAAAAA ATCCCTCCCT CATCTCCCTG ATGAAATCT	9840
CGCGTTGACC ACGTCCACGA TAAAGCTGAA ACTGGTCTTG GCTTGTTCGG GTACCGA	9897

(2) INFORMATION FOR SEQ ID NO: 11:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 8148 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 11:

CCGTGGAAACA AGCCAAGACC AGTTTCAGCT TTATCGTGGG CGTGGTCAAG CCGAATTTT	60
CATCAAGGAG ATGAAGGAGG GATTTTTTGG CGATAAAACG GATAGTTCAA CCTTAATCAA	120
AAACGAAGTT CGTATGATGA TGAGCTGTAT CGCCTACAAT CTCATCTCTT TTCTCAACAA	180
TCTAGCTGGA GGTGACTTCC AAACTTTAAC AATCAAACGC TTCCGCCATC TTTTCTCTCA	240
CGTGTAGGA AAATCTGTTT GAACAGGAGC CAAGCACTC CTCAAATTGT CTAGTCTCTA	300
TGCTATTCC GAATGTTTT CAGCACTTTA TTCTAGGATT AGAAAAGTCA ACCTGAATCT	360
TCCTGTCTCT TATGAACCAAC CTAGAAGAAA AGCGTCGTTA ATGATGCATT AAAGAACAGT	420
CGAGATGAAA AAATCTGTG ACGAACCAAG GGAGGATCT GCCCTTTTGA GGAATCTAG	480
CGAGGAAAAA CGATACTGGA ACGCAGAAAA GTAAAACTGA CCTCATGAGG AGGAAGAAAG	540
TGCTCATGA GGTGAGGGT TTTGTAAGTT ACATCTAGTT GAGAGAGGTA TGAATGATT	600
GGGTAATAC AATGAGCTTG AAAGAAGTAG CAAACTCACC AAGCGCCAAT TCTTTGAGAA	660
TCAGATGCTG GATTATACCA TCATTGCGCA TGAGAGTTT GAAATCATCC GTCATTCTGT	720
CTACCAGACA GATGATCTG AAGTGGAAAA TGCTCTGCT TTTGAAGTGA AAAATGATGA	780
AACAGACAAG CTGATTCTGT TATTAAGCGA GGAATTGGT GTAGGTGAAA AATTGTGCTT	840
CGTTGACGGA ACAAATGTC GTGAAAAATG TTTATATAT GATAAATAA ATGAGAGAA	900
GATTCGCTTG CAGTCTAGA AATAGGCATT TTGAATAGTG AATATGTTAT AATAAGTATT	960
AGTAGAGGT GTTTGAGATT GGAGAGGAAA CTGACCATAA AGACATTGC GGAAATGGCT	1020
CAGACCTCGA AAACAACCGT GTCACTTTAC CTAACGCGGA AATATGAAA AATGTCCCAA	1080
GAGACACGTG AAAAAGTTGA AAAAGTTATT CATGAAACRA ATTACAAACC GAGCATTTGT	1140
GCGCGTAGCT TAAACTCCAA ACGAACAAAA TTAATCGGTG TTTGATTTGG TGATATTACC	1200
AACAGTTTCT CAAACCAAA TGTTAAGGGA ATTGAGGATA TCGCCAGCCA GAATGGCTAC	1260

CAGGTAATGA TAGGAAATAG TAATTACAGC CAAGAGAGTG AGGACCGGTA TATTGAAAGC	1320
ATGCTTCTCT TGGGAGTAGA CGGCTTTATT ATTACAGCGA CCTCTAATTT CCGAAATAT	1380
TCTCGTATCA TCGATGAGAA AAAGAAGAAA ATGGTCTTTT TTGATAGTCA GCTCTATGAA	1440
CACCGGACTA GCTGGGTAA AACCAATAAC TATGATGCCG TTATGACAT GACCCAGTCC	1500
TGTATCGAAA AAGGTTATGA ACATTTTCTC TTGATTACAG CGGATACGAG TCGTTTGAGT	1560
ACTCGAGATTG ACGGGCAAG TGGTTTGTG GATGCTTAA CAGATGCTAA TATGCGTAC	1620
GCCAGTCTAA CCATTGAAGA TAAGCATACG AATTTOGAAC AAATTAAGGA ATTTTACAA	1680
AAAGAAATCG ATCCCGATGA AAAAATCTG GTATTATCC CTAAGTGTG GGCCCTACCT	1740
CTAGTCTTAA CCGTATCAA AGAGTTGAAT TATAACTTGC CACAAGTTGG GTTGATTGTT	1800
TTTGACAATA CGGAGTGGAC TTGCTTTTCT TCTCCAAGTG TTTCGACGCT GGTTCAGCCC	1860
TCTTTTGAGG AAGGACAACA GGCTACAAAG ATTTTGATTG ACCAGATTGA AGGTCCGAAT	1920
CAAGAAGAAA GGCAACAAGT CTGGAATTGT AGTGTGAAT GGAAGAGTTC GACTTCTTAA	1980
AATGAAGGAA AATGACTTGC AATCTCTGTT AAGAAATAAA ATAAATCCAC CTAGAACAAG	2040
CTAGTGGGA TTATTTCCTT ATGAATGAG AAATTATGGG AGCAAGTCC TAAATCAACT	2100
GTTTTTGATC TACTTCTTAA ACTACTTGAT AAAAGTTATA GAAGTAGGCC AAATCTGAAA	2160
TGATGGTTAC GACTAGGAAT ATTGAAAAAT TCCATTGGAC AGGGTTGGTT AAAAGTTGTG	2220
GAAAGGATAT GAGGAGAAAG AAGAGGGCTG CGTTGAGGAC AGGTATCCGT TTTGATTGTA	2280
TTTTCTCAAG TCCTTTATGT AGCCGAGGAA GAAAGAGGAG TAGGAGTAGT AAAATCTAT	2340
GAGAAATAGC TCCGAAATGA AGGCCGAAAG AAAGGAAAAT ACTGATATAA ACATGAATGA	2400
TCAGTAGTCT AGCTAGTAT TTCTATAAGC ACCTCCTAAT CCTGGTCTTT TTTAGTCTCT	2460
GCAATACGAA GTGAGTCGAC AATATGTATC ATCACTCCGA AAAAGAAAGC TCCAGTATA	2520
GTTTTTAAJA TATGTTTTGT ATTTAGAAGA GAATGATAA AATTGGATT TTCACTTGT	2580
AGGGTATCAA TGAGTGGAAT TATAAAAAAT ATCACTGTTC CATAAATCGA ACCTGCTTTC	2640
AGACCAGGAT AACGTAACGT TTTCTTTTCT TTTTTCATGA GTTTCCTCCT AATCCTCATC	2700
TTGATTTTTT TTAGTTTTTG CAATGCGACG GGAGATGAGG AACTGTATGC TCGCTCGGAA	2760
GAAATAGAAA CCGAGATATC TTGATACACC ATTTCTTATA GTGAGAAGAG AATGAAATTA	2820
GTCTGACCT TCATCTATGA GTATCCTGAG AAGAGGAGTT ATAAAAACA TCCATAGACC	2880
AAAGAACAAA CCTGCTTTCA GACCTGGGTA GTGTAGTTGC TTGCTTTCTT TCTCATTCAG	2940
CATATCTGGT TCAATGACTG TGATGCTGT TTTTTCATT TGGTAGGTGA CATAGCCAGA	3000

222	
AGCGATGAGG GCAATCACTA AAATCAGAG AGGATAGATT AGAGCCACTT CTTGAGGCTA	3060
TTTATAGGCC AGAAGCAGTG GAATAAGATT TCGGAAATC ATCAGATAAA AGAGGATGAT	3120
AAAGACTTGG TTCCCAATAC TATCGGCCTC ACGCCGTTG TATTGCTCAA GGGGACCAGA	3180
AATACCGTAT GTGCGTTTGA TCAGTTTTTC AGTGAAGGTT TCTTTTTTCA TGAGTTTGCT	3240
CTTTTTTAA AAATCTTCCT CCCAAAGAG ACTGTTGAG TCAGTTTGA GCGTGGGGC	3300
GAGATTGAGA CAGAGTTCCA AGGTTGGAAT GTACTTGTCTG TTTTCAATCA TATTGATAGT	3360
CTGTCTGAG ACACCGATAT CCTTGGCGAG TTCGAGCTGG GAAATACCCA ATTCTTTCG	3420
AAATCTTTC ACACGATTCA TCTGTTCTCC TTTCTGATTT ATGTCGTATA TATTTGACTA	3480
TATTATAGTC TTTTAAACAT AAAGTGTCAA GTATTTTTGA CATATTTTT GAAGAAATAG	3540
TAGTCTCCTT GTCTTATTTG TCTGACAAGT GCAAGCTGGT CGGATTTCTG GTAAATAGA	3600
TAAGATATGA CAAAAGAATT TCATCAATGA ACGGCTTAC TCCACGAAAC GATTGATATG	3660
CTTGAGCTAA AGCCTGATGG TATCTACGTT GATGCGACTT TGGGCGGAGC AGGACATAGC	3720
GAGTATTTAT TAAGTAAATT AAGTGAAGAA GSCCATCTCT ATGCTTTTGA CCAGGATCAG	3780
AATGCCATTG ACAATGCGCA AAAACGCTTG GCACCTTACA TTGAGAAAGG AATGGTGACC	3840
TTTATCAAGG ACAACTTCCG TCATTTACAG GCATGTTTGC GCGAAGCTGG TGTTCAGGAA	3900
ATTGATGGAA TTTGTTATGA CTTGGGAGTG TCTAGTCTCT AATTAGACCA GCGTGAGCGT	3960
GGTTTTCTT ATAJAAAGGA TCGGCCACTG GACMTGCGGA TGAATCAGGA TGCTAGCCTG	4020
ACAGCCTATG AAGTGGTGAA CAATTATGAC TATCATGACT TGGTTCGTAT TTTCTTCAAG	4080
TATGGAGAGG ACAAAATCTC TAAACAGATT GCGGTAAAG TTGAGCAAGC GCGTGAAGTG	4140
AAGCCGATGT AGACAACGAC TGAGTTAGCA GAGATTATCA AGTTGGTCAA ACCTGCCAAG	4200
GAACTCAAGA AGAAGGGGCA TCCTGCTAAG CAGATTTTCC AGGCTATTGC AATTGAAGTC	4260
AATGATGAAC TGGGAGCGGC AGATGAGTCC ATCCAGCAGG CFATGGATAT GTTGGCTCTG	4320
GATGGTAGAA TTTCACTGAT TACCTTTTAT TCCTTAGAAG ACCGCTTGAC CAAGCAATCG	4380
TTCAAGGAAG CTTCAACAGT TGAAGTTCCA AAAGGCTTGC CTTTCATCCC AGATGATCTC	4440
AAGCCCAAGA TGGAAATGGT GTCCCGTAAG CCAATCTTGC CAAGTGCGGA AGAGTAGAA	4500
GCCAAATACC GCTCGCACTC AGCCAAAGTT GCGTGGTGA GAAAAATCA CAAGTAAGAG	4560
GGAAAAAGAT GGCAGAAAAA ATGGAAAAAA CAGGTCAAAT ACTACAGATG CAACTTAAAC	4620
GGTTTTCGCG TGTGAAAAAA GCTTTTTACT TTTCCATTGC TGTAAACCACT CTTATTGTAG	4680
CCATTAAGTAT TATTTTTATG CAGACCAAGC TCTTGCAGT GCAGAAATGAT TTGCAAAAA	4740
TCATATGCCA GATAGAGGAA AAGAAGACCG AATTGGACGA TGCCAAAGCAA GAGGTCAATG	4800

AACATTACG	TGCAGAACGT	TTGAAAGAAA	TTGCCAATTC	ACACGATTGG	CAATTAAJCA	4860
ATGAAAATAT	TAGAATAGCG	GAGTAAGATA	TGAAGTGGAC	AAAAAGAGTA	ATCCGTTATG	4920
CGACCAAAAA	TCGGAAATCG	CCGGCTGAAA	ACAGAOCGAC	AGTTGGAAAA	AGTCTGAGTT	4980
TATTATCTGT	CTTTGTTTTT	GCCATTTTTT	TAGTCAATTT	TGCGGTCATT	ATTGGGACAG	5040
GCATTCGCTT	TGGAACAGAT	TTAGCGAAGG	AAGCTAAGAA	GGTTCATCAA	ACCACCCGTA	5100
CAGTTCCTGC	CAACGCTGGG	ACTATTATAT	ACCGAAATGG	AGTCCCGATT	GCTGAGGATG	5160
CAACCTCCTA	TAATGTCTAT	GCGGTCATTG	ATGAGAACTA	TAAGTCAGCA	ACGGGTAAGA	5220
TTCTTTACGT	AGAAAAACA	CAATTTAACA	AGGTTCGAGA	GGTCTTTCAT	AAGTATCTGG	5280
ACATGGAAAG	ATCCTATGTA	AGAGAGCAAC	TCTCGCAACC	TAATCTCAAG	CAAGTTTCTT	5340
TTGGAGCAAA	GGGAAATGGG	ATTACCTATG	CCAATATGAT	GTCTATCAAA	AAAGAAATGG	5400
AAGCTGCAGA	GGTCAAGGGG	ATTGATTTTA	CAACCAAGTC	CAATCGTAGT	TACCCAAACG	5460
GACAAATTGC	TTCTAGTTTT	ATCGGCTCTAG	CTCAGCTCCA	TGAAATGAA	GATGGAAGCA	5520
AGAGCTTGCT	GGGAACCTCT	GGAATGGAGA	GTTCCTTGAA	CAGTATTCTT	GCAGGGACAG	5580
ACGGCATTAT	TACCTATGAA	AAGGATCCTC	TGGGTAATAT	TGTACCCGGA	ACAGAAACAAG	5640
TTTCCCAACG	AACGATGAC	GSTAAGGATG	TTTATACAAC	CATTTCAGC	CCCCCTCAGT	5700
CCTTTATGTA	AACCCAGATG	GATGCTTTTC	AAGAGAAGGT	AAAGGGAAG	TACATGACAG	5760
CGACTTTGCT	CAGTGCTAAA	ACAGGGGAAA	TTCTGGCAAC	AACGCAACGA	CCGACCTTTG	5820
ATGCAGATAC	AAAAGAAGGC	ATTACAGAGG	ACTTTGTTTG	GCGTGATATC	CTTTACCAAA	5880
GTAACATGTA	GCCAGGTTCC	ACTATGAAAG	TGATGATGTT	GGCTGCTGCT	ATTGATATAA	5940
ATACCTTTCC	AGGAGGAGAA	GTCTTTAATA	GTAGTGAGTT	AAAAATTGCA	GATGCCACGA	6000
TTTCGAGATTG	GGACGTTAAT	GAAGGATTGA	CTGGTGGCAG	AACGATGACT	TTTTCTCAAG	6060
GTTTTGCACA	CTCAAGTAAC	GTGGGATGA	CCCTCCTTGA	GCAAAAGATG	GGAGATGCTA	6120
CCTGGCTCTGA	TTATCTTAAT	CGTTTAAAT	TTGGAGTTCC	GACCCGTTTC	GCTTTGACGG	6180
ATGAGTATGC	TGGTCAGCTT	CCTGCGGATA	ATATTGTCAA	CATTGGCGAA	AGCTCATTTG	6240
GACAAAGGAT	TTCACTGACC	CAGACGCAAA	TGATTCGTGC	CTTTACAGCT	ATTGCTAATG	6300
ACGGTGTCAAT	GCTGGAGCCT	AAATTTATTA	GTGCCATTTA	TGATCCAAAT	GATCAAACTG	6360
CTCGGAATTC	TCAAAAAGAA	ATTGTGGGAA	ATCCTGTTTC	TAAAGATGCA	GCTAGTCTAA	6420
CTCGGACTAA	CATGGTTTTG	GTAGGGACGG	ATCCGGTTTA	TGGAACCATG	TATAACCACA	6480
GCACAGGCAA	GCCAACTGTA	ACTGTTCTCT	GGCAAAATGT	AGCCCTCAAG	TCTGGTACGG	6540

224

CTCAGATTGC TGACGAGAAA AATGGTGGTT ATCTAGTCGG GTTAACCGAC TATATTTTCT	6600
CGGCTGTATC GATGAGTCCG GCTGAAAATC CTGATTTTAT CTGTGTATGT ACGGTCCAAC	6660
AACTTGAACA TTATTAGCTT ATTCAGTTGG GAGAAATTCG CAATCCTATC TTGGAGCGGG	6720
CTTCAGCTAT GAAAGACTCT CTCATCTCTC AAACAACAGC TAAGGCTTTA GAGCAAGTAA	6780
GTCAACAAGG TCCCTTATCT ATGCCTAGTG TCAAGGATAT TTCACTTGTG GATTTAGCAG	6840
AAGAAATTCG TCGCAATCTT GTACAAACCA TCGTTGTGGG AACAGGAACG AAGATTTAAA	6900
ACAGTTCTGC TGAAGAAGGG AAGAATCTTG CCCCAGAACCA GCAAGTCCCT ATCTTATCTG	6960
ATAAAGCAGA GAGAGTTCCA GATATGTATG GTTGGACAAA GGAGACTGCT GAGACCTTG	7020
CTAAGTGGCT CAATATAGAA CTGAAATTC AAGGTTCTGG CTCTACTGTG CAGAAGCAAG	7080
ATGTTCTGTC TAACACAGCT ATCAAGGACA TTAATAAAT TACATTAAC TTAGGAGACT	7140
AATATGTTTA TTTCATCAG TGCTGGAAT GTGACATTT TACTAACTTT AGTAGAAAT	7200
CCGCCCTTTA TCCAAATTTA TAGAAGGCGC CAAATTACAG GCCAGCAGAT GCATGAGGAT	7260
GTCAAAACAC ATCAGGCAAA AGCTGGGACT CCTACAAATG GAGGTTTGGT TTCTCTGATT	7320
ACTTCTGTTT TGGTTGCTTT CTTTTGCGC CTATTTAGTA GCCAATTCAG CAATAATGTG	7380
GGAATGATTT TGTTCACTTT GGTCTGTAT GGCTTGGTCG GATTTTAGA TGACTTCTC	7440
AAGGTCTTTC GTAAATCAAA TGAGGGGCTT AATCCTAAGC AAAAATTAGC TCTTCAGCTT	7500
CTAGGTGGAG TTATCTCTTA TCTTTCTAT GAGCGCGGTG GCGATATCCT GTCTGCTTT	7560
GGTTATCCAG TCAATTTGGG ATTTTCTAT ATTTCTTCTG CTCTTTCTG GCTAGTCGT	7620
TTTTCAACAG CAGTAAACTT GACAGACGGT GTTGACGGTT TAGCTAGTAT TTCCGTTGTG	7680
ATTAGTTTGT CTGCCTATGG AGTTATTGCC TATGTGCAAG GTCAGATGGA TATCTTCTA	7740
GTGATTTCTG CCATGATTGG TGGTTTGCTC GGTTCCTTCA TCTTTAACA TAAGCCTGCC	7800
AAGGTCTTTA TGGGTGATGT GGAAGTTTG GCCCTAGGTG GGAATCTGC AGCTATCTCT	7860
ATGGCTCTCC ACCAAGAATG GACTCTCTTG ATATCTCGAA TTGTGTATGT TTTTGAACA	7920
ACTTCTGTTA TGATGCAAGT CAGTTATTTT AAATGACAG GTGGTAAACG TATTTCCGT	7980
ATGACGCCCTG TACATACCA TTTTGAGCTT GGGGATTGT CTGGTAAAGG AATCTCTTG	8040
AGCAGTGGGA AGGTGACTT CTCTTTTGG GAGTGGGAC TTCTAGCAAG TCTCCTGACC	8100
CTAGCAATTT TATATTTGAT GTAAGAATGG CACCTGATG TTTTCAGGG	8148

(2) INFORMATION FOR SEQ ID NO: 12:

- (i) SEQUENCE CHARACTERISTICS:
  - (A) LENGTH: 9909 base pairs
  - (B) TYPE: nucleic acid

225

(C) STRANDEDNESS: double  
(D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 12:

TACTCCACC TTAATATCCG TTCTGTAAA TACTTTACCG CTTTAAAGTT CATAGAATTG	60
AACCTTTTAA TGCTGTCTT CAAGCATCTT TTCCATCCAA TTTTTAGGAG TTTGACCAGC	120
TTTAAATAAA AACCTTGCTG GGTGTAGTAG TATAGATTTA TCTCGATTT TATAAGCTTC	180
ATCAATAAAA TAGTGATATA TCGGCTCAT TCTGGCTTCT CCTGTTCCT GATACGGAGG	240
ATTTCTATC ACGACATCAA ATTTCAATTC ACTTCTCTG CTAGATAGGC GCTCAAAACC	300
TATCATTTTA TTCTTTTCC AGTCTTGAT ATGCTTTTA GATTCTCTA CTCTGTGAC	360
TTCTAGCTCA TCCGCAACA AACTCAATTG TTGAGATTGC TTTTGTTAG CTGAATAAGG	420
ACTACTTTT TTCAATCCAT CCATCTGAAA GACATTGTAA GAGATAATAG TCGCAATTC	480
TTTCTTTTGC TCTAATGTTG GTTGATTTCC AGTCTTAGCT AGATAATAGT CCTCAAAAGT	540
TGCCAAAGA TTCTCACGCG CCAAAAGGAG AGAATCTCT TGATACTCAT AACCATACGA	600
AGCATGATAA GCATCTTTTA CAAGTTTATA AATGTGACT TCATCTGAAA CCTCACGACT	660
AATCCGTGTC AGTTTCTAT CAACAAAACC AACTCGCTCA GATAATGGAA TTTCTTCACC	720
AGTTACGGTA TCATATCTCG TTACCATATA AGTGTCTTA CCACAAGTAA CTTCTAACCA	780
TCGTAAAGTC ACATACTCT CAAGACTTAA CGAGCCTAAT TTGATTCTA CATATCCATT	840
TTGCTTTGCG ACCAACCAAG TTGGTGTAAA CACTTCTGCC CTTATTTTTG TCCGATCTTT	900
TTGTTCATAT TTGGATTTT CAGATCTGGG CTGAATCAAG TTGGCAAAGT TTCCAGTAAC	960
CTTACTTGA TTGATGOGAT CACTTGGAGC AAATCCCTTT CCTAACAATT CATAAGAATT	1020
CGTAGGCCAA ACAATTGATT TCTTTGCTGT TCGATCTTT AAAAGAATTT TTAAATAAGTC	1080
AGCGGATCT TTAGCCAAAC TTTCTTCACT AATATCTATT GTCATCAGCA ACCTCTCTTA	1140
TATTGTAGC CTAATTATAT CATATTTTAA AGAATGAAA TTTACTTGA AAAAGTAATT	1200
CAATAAATAT CTCTCCGATG ACCAACTTCT AGAGTAGCAA CGACTAATTC ATCATCTACA	1260
ATTTGTACGA TAACTCGATA ATTACCAATT CTATAGCGCC ATTGACCAAC GCGATTACCA	1320
ACCAAGCCT TTCCGTGTCG TCTTGGGTCT TCCAAAACAT TGGTTGTAA ATAGTTGTGA	1380
ATTAGCTTCT GCGTATAACG GTCCAATTTT TTCAATTGCT TGATAAAACG TCTTGTGGA	1440
ACTAATTTAT ACAAAATTAT CATCTTCAA GCTAATACT TGATCATATT CTTCOCAAGT	1500
AATGGGTCA ACTCCTTTT CCAAGTCTC TAAATACTCT TGATAGGCTA AATCTGCCAC	1560

		226	
ACGAGCATCG	TATTCATCTT	CFAGGGCTTC	AAGAGTTTTC
GTGCGAATPA	GTTCGAAAG	1620	
GGAAACTCCT	TCAAACCTAG	CCATTCCTTT	CATAAATGTT
TTATCAGCTT	CAGAACTTT	1680	
TAATGTATA	GTAGTCATCT	TTTGTCCTCC	CTTTTFFAAT
GGTAACACCA	TTCATTACT	1740	
TTTTAGGTGT	TCAGTCAATA	TAAAAAGAAC	ACCTTCCTCAG
CGTTCCTTCT	ATATCTCTGT	1800	
CAATGGTGT	GCGGTATCTG	GTGAGGTATC	ATAAACCTTA
AAGTCTACTC	CGACTCCAG	1860	
ATCAGCTTGA	GCCAGCTGAT	TGACCATGGT	CATATGAGCC
AGTTCCTTGA	TATGTGTTTC	1920	
CTTAGATAAA	TGCCCAAGGT	AAATCTTCTT	AGTACGATTT
CCTAGCGTCC	GAATCATAGC	1980	
TTCAGCACCG	TCTCTGTTAG	AAAGGTGACC	AAGGTGAGAT
AGGATTCGTT	GTTCGAGTCG	2040	
CCAAGCGTAA	GAACCTGATC	GCAAAATCTC	TACATCATGG
TTGGCCTCGA	TAAGATAACC	2100	
ATCCGCAATT	TCGACAATGC	CCGCCATAGC	GTCACTGACA
TAACTCTGAT	CTGTCAAGAG	2160	
GACAAAATCT	TTATCATCTT	TCATAAAGCG	ATAGAAGCTC
GTCGCGACTG	CATCATGGCT	2220	
TACACAAAA	CTCTCGATGT	CGATATCTCC	AAAGGTTTTG
GTTCCTTCCA	TTTCAAAAT	2280	
ATGCTTTTGC	GAAGAATCCA	CCTTGCCAAG	ATATTTACTA
TTTTCCATAG	CTTGCCAGGT	2340	
CTTTTCATTG	GCATAAAGAT	CCATACCAT	CTTGCGAGCC
AAAACGCTCA	CTCCATGGAT	2400	
ATGATCTGAA	TGCTCATGGG	TAAATCAAGAT	GGCATCCAGG
TCTTCTGGCT	TAAGGTTAAT	2460	
TTTCAGCTAGC	AGACTGGTAA	TTTTCTTGCC	AGACAAGCTT
GCATCTACTA	AAAGCTTCTT	2520	
TTTTGAGGTT	TCCAGATAAA	AAGAAATTC	ACTGGAACCC
GACGCTAAAA	TACTGTATTT	2580	
AAAGCCTATT	TCACCTCATC	TAGTCTTCTA	CTTCATCTCT
CCATCTCTCT	TCTTTCACCTG	2640	
CATCCTTATC	ATAAGGGAGT	ACAATGGTAA	AGGTTGAACC
CTTGCCGTAT	TCACCTCTTG	2700	
CCCAATATAA	GCCCTTATCT	TGTTTGATAA	TTTCTTTAGC
GATAGACAGT	CCTAGACCTG	2760	
TACCACCTTG	TGCACGACTT	CTAGCACGAT	CCACACGATA
GAAACGGTCA	AAGATACGTG	2820	
GTAAATCCTG	CTTAGGAATC	CCCAAAACCGT	GGTCAGAAAT
GGATAAAATC	ATCTGGTCTT	2880	
CAGTTGTCTT	CATTCTGACA	GTGATTTTAC	CCCCATCTGG
CGAATACTTA	ATAGCATTAT	2940	
TTAAAAATAT	GTGACAACC	TGCGTCATCT	TATCTGTATC
AATTTCATC	CAGATAGAAAT	3000	
TGATGGGATA	ATCTCTCACC	AACTCATATT	TTTTCTCCTT
TCTCTGCTCT	TTTCTCTCTT	3060	
CAAAACGATT	GAGGATAAAG	GTAAATAAAG	CAGTGAAGTT
AATCAGTTCC	ACATCTAGTC	3120	
GACTGGTAGC	ATTATCAATA	CCTGAAAGAT	GGAGGAGATC
CGTCACCATG	CGCATCATAC	3180	
GTTTGGTCTC	ATCAAGAGAA	ACCTTGATAA	AGTCTGGTGC
TACAGTTTCA	CACAAAGCCC	3240	
CCTCATCCAA	GGCTTCAAGA	TAGGATTTTA	CGCTAGTCAG
AGGAGTCCGT	AACTCATGGC	3300	
TAAACATTGGA	AACAAGAGAT	CTTCGTTGCG	GTTCCTTCCTT
CTCCTGCTCC	GTGCTATCAT	3360	



GCAAAACAGC CACCAACCT GAAATAAAGC CAGACTCTCG ACOTATCAAG GCAAAGCGAA	3420
CTCGAAGGTT CAAATATTTCG CCATTGATAT CTGGGAATC TAGCAACAAT TCTGGACTTT	3480
GGGTAAATCAA ATCAAGCAAT TCATAGITTT CTCTATCTT GAGCAATTCC AAAATGCTTC	3540
TATTCAGAAC ATCTCTCTTA ACCAACCCCA GTTGTCTCTT GGCTGTATCG TTAATCATGA	3600
TAACTGTACC CCGACGGTTA GTCGCAAGAA CCCCATCTGT CATATAAAAC AGAATACTAT	3660
TTAGCCTCTT ACTCTCTTGT TCTAGATTTT CCGAGTGAG ACGAATAACC TCCGACAAAT	3720
CATTCAAATT ATTGGTAATA TTGGTGATTT CAGACCCACC TTGCATATCA AGAACCCTGG	3780
AATAATCTCC TGCAATCAAA TCTTTAACC TTTGATTGAC TTGCTTCAAC TGAATATTAT	3840
CACGCTATT TTCCAGTAAT AAGAGGGTCA CAACAAGGAT GAAACCTAAC AAAATCAGGA	3900
TAAAGATAAA ATCTCTGGTA AAAATGGTTT GTTTCAGTAA ATCAAGCATT ATTTCTCATG	3960
TAAATCCCTA CACCACGGCG CGTCAAGATA TACTCTGGTC GGCTGGCGGT ATCTTCAATC	4020
TTCTCACGCA GACGTCGTAC AGTCACATCA ACTGTACGGA CATCACCAA ATAGTCATAA	4080
CCCCAGACAG TCTCAAGCAA GTGTTGCGGC GTGATGACTT GACCTGTATG CGATGCTAAA	4140
TGATACAAAA GCTCAAAATC ACGATGGGTT AAGTCTAGTT CTTCGCCATA TTTTPTAGCC	4200
ACGTAGCGGT CTGGAACAAT TTCTAAATCC CCAATTGGGA TAGGTGAGG TTTACTATAT	4260
GCTTCTCTGAC CATCTACTGG CATAGGTTGA GAACGACGCA GAAGAGCTTT AACACGCGCC	4320
TGCAACTCAC GATTGGAGAA GGGTTTGT TTACATGATC CTGCCCAAG TTCCAAACCG	4380
ATAACCTTAT CAAATTCACT ATCTTTGGCT GAAAGCATAA GAATGGGCAC ACTGCTTGTG	4440
TTACGAATGG TCTTAGCAAC TTCTAAACCA TCAATTTCTG GAAGCATCAA ATCCAGAATA	4500
ATAATATCTG GTTGTCTCTG TTCAAAATGC TCTAGCGCTT CACGACCATT AAAAGCAGTT	4560
ACAACTTCGT AACCTTCTTT GGTCAATATTA AACTTGATTA TATCCGAGAT TGGTTTCTCA	4620
TCATCTACAA TTAGTATTTT TTTCATATGT TCACCTTTTT CTCTACTATP ATACCAAAAA	4680
AATAGTCAGA AGACACAATA GCTAGTCTTG GCTACTGTCT AAGTTGGCTT GTGCATAAAC	4740
CTGCCAGATT TTTTGTGGG GTTTGGCAAG TGGGTAATTC TTGAATCTT CTGTGAAAG	4800
CCAGCGAACT TCCCTATCTG AAAAATCATG GAAOTCACTC ACCTGACCTG CTACAATCTG	4860
TACATGCCAT TTTTCATGAC TAAAAACATG CTGGACTGTA TCAAAACAAA CATCAAGCCA	4920
ATCAACATCT AGGTCAATAG CCGTCTGGAA ACTCTCTTCT GGACTGGGAC CAAAGTTTCA	4980
ACTTCTTCC GCAACCTGAT GAAAGAGGTC AACTGCTCTT TCTTGGGAAA AGTTATCAAC	5040
TTCTATAAAG GGGAAATGCC AAAACCTGTC CAAGAGCTTT TGCTTTTCAAT TTTTTCAAAG	5100

228	
TAAAAATGTT CCTTGAGAA TTTTCACAAC TAAGGCTTPTA AGATAAATPAG GAACCGGCTT	5160
TTTCTTTAGGA GATTTAATTG GATAACGGTC CATGGTTCCA TTCTGATATG CCGCACTAAA	5220
GTCTTGTGACT GGGCTTTCTT CAGGTCCTGGG ATTTACAGGA GACTCAATAT CAGACCTTAA	5280
GTCCATCAAG GCTTGTATTA AATCACCCTG ACGATCCGGA TTAATCAAGA TCTCCATCAT	5340
TGCCCTGAAAA ATTTTTCGAT TACTTGGAA TCCCAATATG TGGTTGACTT CAAACAGACG	5400
CGCCAGAACG CGCATGACAT TACCATCTAC AGCTGGCTCA GGCAAGTTAA AAGCAATACT	5460
GGAAATGGCT CCTGCTGTGT AAGGTCCAAT CCCTTTCAGG CTGGAAATTC CTTCATAGGT	5520
ATTTGGAAAT TGGCCACCAA AGTCAGTCAT AATCTGCTGG GCTGCAGCCT GCATATTGGG	5580
AACTCGAGAA TAATAGCCCA AGCCCTCCCA AGCTTTCAGT AAATCTCTCT CAGGCGCAGT	5640
TGCCAGACTT TCGACAGTGT GAAACCAGTC CAAAAATCTT TCGTAGTAAG GGATAACTGT	5700
ATCCACCCTG GTCTGCTGAA GCATGATTTT AGATACCGAG ATGTGATAAG GATTTTCTT	5760
TCTTCCTCAA GGCAAATCTC TTTTGTTC ATCATAACAA CGGAGAAAGT TCTCACGGAA	5820
AGAAATGACT TTCTCTCCG GGCACATGAC GATACCGTAT TCTTTCAAAT CTAACATATC	5880
TCTAGTATAA CACAGAAGST TTCACCTGTC TTGTATCTGT ATTTATAATA TTTTCAATAG	5940
ATAGTATATA ACTTTTCTAT CTACTTATAC TCAATGAAAA TCAAAGAGCA AACTAGGAAG	6000
CTAGCCGACG GTTGCTCAA ACACGTGTTT GAGGTTGTGG ATAGAAGTGA CAGAGTCAGT	6060
ATCATATACT ACGGCAAGST GAAGCTGACG TAGTTTGAAG AGATTTTCTGA AGAGTATAAA	6120
TCTTATTGAT GAACTGCTTG CAGTCTGAGA AAAAAAGAGC TTGGATATTA TTTCCAAACT	6180
CACTTAAAGT CAATTTCAAT CCACTAGAAC AAGCTAGTA CAGTTCCATC GCTTTCACCA	6240
TCCATGTGGA GAGCTGCTGG ACGTTTGGGA AGACCTGGCA TGGTCATAAC ATCAACAGTT	6300
AAGGCAACGA TGAAGCCTGC ACCTAATTTT GGTACCAATT CACGAATGGT AATTTCAAAG	6360
TTTTCCTGGG CTCCAAGCCG ATTTGGATTG TCTGAGAAAC TGTATTGAGT TTTAGCCATA	6420
CAGATTGGCA ATTTGTCCCA ACCGTTTGA ACGATTGAG CAATTTGTGT TTGAGCTTTC	6480
TTCTCAAAGT TCACTTTGCT ACCACGATAG ATTTCACTGA CAATTTTTC AATCTTTCT	6540
TGGACAGAAA GGTCAATTATC ATACAAACGT TTATAGTAGT CTGGATTTTC AGCAATTGTC	6600
TTAACAAAGT TTTGGGCAAG TGCTACTCCA CCTTCTGCTC CATCAGCCCA GACACTAGCT	6660
AATTCACCTG GTACATCGAT TGAAGGCACG AGTCTTTTAA AGGCTCAAT TTCACTTCT	6720
GTATCAGATA CAAATTCGTT AATAGCTACA ACTGCTGGAA TACCGAACTT ACGGATATTT	6780
TCAACCTGGC GTTTCAGGTT AGCAAAACCT GCACGAAGTC CCTCTACATT TTCTTCAGTC	6840
AGAGCGTCTT TAGGCCACCC ACCATTTCATC TTAAGGGCAC GAAGGGTTGC GACAATAACA	6900

ACTGCATCTG GAGATGTTGG CAAGTTTGGT GTCCTGATAT CAAGGAATTT CTCAGCACCA	6960
AGGTCCGCAC CAAAACAGC TTCAGTAACA GTGTAATCAG CCAAGTGAAG GGTGTTTGT	7020
GTGCCCAAA CAGAGTTACA GCCATGAGCG ATATTGGCAA ATGGACCACC GTGTACAAAG	7080
GCAAGTGATC CGTAAATTGT CTGAACCAAG TTTCGGCTTAA TAGCATCCTT CAAAATCAAA	7140
GCCAAAGCAC CCTCAACCTG CAAATCACCT ACAGAAACAG CGGTACGGTC ATAGCATATA	7200
CCAATAACGA TATTCGCCAA ACGACGTTTC AAGTCCTCGA TGTCCTTGC CAAGCAAAGA	7260
ATTGCCATGA TTTCTGAAGC AACTGTAATA TCAAAACCAT CCTCAGGTGG AATACGGTTT	7320
AGAGGACCAC CAAGACCAAC AGTCACATGG CGGAGCGTAC GGTCTGTCAA GTCCACAAAG	7380
CGTTTCCAGA GGATACGACG TTGATCAAIT CCAGCTCAT TCCCTTGGTG CAAGTGGTTG	7440
TCAATCAAGG CAGAAAGGGC ATTGTTGCCA GTTGTAAATAG CATGCATATC TCCAGTAAAG	7500
TGGAGGTTGA TGCTCTCCAT TGGCAGAACT TGTGCATACC CACCACCAGC AGCACACCC	7560
TTGATCCCCA TGACTGGACC AAGAGACGGT TCOCGGATAG CAATCATGGT TTTCTTGCCA	7620
ATCTTGTTC AAGCATCCGC AAGACCAATG GTAAAGCGTC ACTTTCCTTC ACCTGCAGGT	7680
GTTGGGTTGA TGGCAGTAAC CAAGATCAAT TTACCGACTG GATTGCTCTC AACTGCACGA	7740
ATTTTATCAA AGCTGAGTTT AGCCTTGTAC TTTCCGTACA ACTCCAAATC GTCAATAAGAA	7800
ATACCAAGTT TCTCTAACAC ATCAACAATT GGCTTCAACT CAATACTCTG TGGGATTTC A	7860
ATATCTGTTT TCATTCAAAA TTTCTCTAAC CTCCTATATG ATAATTCATT ATATCAGAAA	7920
ACAAGATTTT TAACATCCTA AAACCTCTCTA AACGTTGTA AATATCTCTG TTTTAAAGAC	7980
TTTTAGAGTC CTTTCTTAAA TTTTATATGG CTTTATAGTT TGAAACTATA ATAAATCTTC	8040
GTTFATACCA AAAATTTATC ACTTTCAATT TACTTACCGC TTATTTTGTG GTACAATAGT	8100
GCTATGAAAA TTTTAGTTAC ATCGGGCGGT ACCAGTGAAG CTATCGATAG CGTCCGCTCT	8160
ATCACTAACC ATTCTACAGG TCACTTGGGG AAAATTATCA CAGAGACTTT GCTTCTGCA	8220
GGGTATGAAG TTTGTTTAA TACGACAAAA CGAGCTCTGA AGCCAGAGCC TCATCCTAAC	8280
CTAAGTATTC GAGAAATTAC CAATACCAAG GACCTTCTAA TAGAAATGCA AGAACGTGTT	8340
CAGGATTATC AGGTCTTGAT CCACTCAATG GCTGTTCTG ACTACACTCC TGTTTATATG	8400
ACAGGGCTTG AGGAAGTTCA GGCTAGCTCC AATCTAAAG AATTTTAAAG CAAGCAAAAT	8460
CATCAGGCCA GATTTCTTTC AACTGATGAG GTTCAGGTTT TGTTCCTTAA AAAGACACCC	8520
AAATCATAT CCCTAGTCAA GGAATGGAAT CCTACTATTC ATCTGATTGG TTTCAAACTG	8580
CTGGTTGATG TTACCGAAGA TCATCTGGTT GACATTGCAC GAAAAAGTCT TATCAAGAA	8640

230

CAAGCAGATT TAATCATCGC GAATGACCTG ACTCAAATTT CAGCAGATCA GCACCGAGCT	8700
ATATTTTGTTG AGAAAAATCA GCTTCAJACA GTCCAGACTA AAGAAGAAAT TGCAGAAATC	8760
CTCCTTGAAA AAATTCAGC CTATCATTTCT TAGAAGGAA AACTATGGCA AACATTCTCT	8820
TGGCTGTAAAC GGGTTCATC GCCCTTATA AGTCGGCAGA TTTAGTCAGT TCTCTAAAAA	8880
AACAAGGCCA TCAAGTCAC TGTCTAATGA CTCAGGCTGC TACAGAGTTT ATCCAACTTT	8940
TGACACTACA GGTACTCTCA CAGAACTCTG TCCACTTGGG TGTCTGAAG GAACCCATTC	9000
CTGATCAGGT CAATCATATC GAACCTGGAA AAAAAGCAGA TTTATTTATC GTGGTAOCTG	9060
CAACTGCTAA CACTATTGCA AACTAGCTC ACGGATTGCG GGACAACATG GTAACAGTA	9120
CAGCTCTAGC CTTACCAAGT CATATTCCCA AACTAATAGC TCTGCTATG AATACAAAAA	9180
TGTATGACCA TCCAGTAACT CAGAAATATC TGAAAAACTT AGAAACTACG GCTATCAGCT	9240
GATTGCTCCT AAGGAATCCC TACTAGCTTG TGGAGACCAC GGACGAGGAG CTTTAGCTGA	9300
CCTCACAAAT ATTTTAGAAA GAATAAAGGA AACTATCGAT GAAAAACGC TCTAATATTG	9360
CACCCATTGC TATCTTTTTT GCTACCATGC TCGTGATACA CTTTCTGAGC TCACCTTATCT	9420
TTAACTTTTT TCCATTTCOA ATCAAAACGA CCATTGTTCA TATTCCTGTC ATTATTGCCA	9480
GCATATTATA TGGTCCACGA GTTGGGGTTA CACTTGATTT TTTGATGGGA TTACTTAGCT	9540
TGACGGTTAA CACGATTACG ATTCTACCGA CAACTACCTT CTTCTCTCCC TTGCTACCAA	9600
ACGGAAACAT CTACTCAGCT ATCAATTGCCA TCGTCCACG TATTTTGATT GGTATAATCT	9660
CTTACTTAGT CTATAAACTG ATGAAJAAAC AGACTGGTCT GATTTTAGCT GGAGCCCTTG	9720
GTTCTGTGAC AAATACTATC TTTGTCCCTG GAGGAATCTT CTTCTATT TTGAAATGTTT	9780
ATAATGGAAT TATCCAACTT CTCTGGCAA CCGTTATCTT AACAAATCA ATGCTGAAAT	9840
TGGTCAATTC TGCAATCTTA ACCCTAGCCA TTGTTCCACG ACTACAAACC TTGAAAAAAT	9900
AAAAACAGG	9909

(2) INFORMATION FOR SEQ ID NO: 13:

## (i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 1126 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 13:

TAATTTTCAT ATAAAGTAA AATAGATGT GTGATTCAAT AATCACCTCA AATGAAAGG	60
AAATTTCTATG TCAAACTCAT CTGTAAATGC AATTCGTTTT CTAGGTATTG ACGCCATTA	120

231

TAAAGCCAAC TCAGGTCATC CAGGTGTGGT TATGGGAGCG GCTCCGATGG CTTACAGCCT	180
CTTTACAAAA CAACCTCATA TCAATCCAGC TCAACCAAAAC TGGATTAAACC GCGACCGCTT	240
TATTCCTTCA GCAGGTCATG GTTCAATGCT CCTTTATGCT CTCTCTCACC TTCTCTGTTT	300
TGAAGATGTC AGCATGGATG AGATTAGAG TTTCCCTCAA TGGGGTTCAA AAACACCAGG	360
TCACCCAGAA TTTGGTCATA CGGAGGGAT TGATGCTACG ACAGGTCCTC TAGGGCAAGG	420
GATTTCAACT GCTACTGGT TTGCCAAGC AGAACCTTC TGGCAGCCA AATATAACCG	480
TGAAGGTTAC AATATCTTTG ACCACTATAC TTACGTTATC TGTGGAGACG GAGACTTGAT	540
GGAGGTGTG TCAAGCGAGG CAGCTTCATA CGCAGGCTTG CAAAACCTTG ATAAGTTGGT	600
TGTTCTTTAT GATTCAAAATG ATATCAACTT GGATGTTGAG ACAAGGATT CCTTTACAGA	660
AAGTCTTCGT GACCGTTACA ATGCTTACGG TTGGCATACT GCCTTGGTTG AAAATGGAAC	720
AGACTTGGAA GCCATCCATG CTGCTATCGA AACAGCAAAA GCTTCAGGCA AGCCATCTTT	780
GATTGAAATG AAGACGGTGA TTGGATACGG TTCTCCAAAC AAACAAGGAA CTAAATGCTGT	840
ACACGGCGCC CCTCTTGGAG CAGATGAAAC TGCATCACT CGTCAAGCCC TCGGTTGGGA	900
CTACGAACCA TTTGAAATTC CAGAACAAAT ATATGCTGAT TTCAAAGAAC ATGTTGCAGA	960
CCGTGGCGCA TCAGCTTATC AAGCTTGGAC TAAATTAGTT GCAGATTATA AAGAAGCTCA	1020
TCCAGAACTG GCTGCAGAAG TAGAAGCCAT CATCGACGGA CGTGATCCAG TCGAAGTGAC	1080
TCCAGCAGAC TTCCAGCTT TAGAAAAATGG TTTTCTCAA GCAACT	1126

## (2) INFORMATION FOR SEQ ID NO: 14:

- (i) SEQUENCE CHARACTERISTICS:  
 (A) LENGTH: 2520 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 14:

CCGGCAACAA AAAAGAAAA ATCAACAGTT AAAAAAATC TAGTCATCGT GGAGTCGCTT	60
GCTAAGCCAA GACGATTGAA AAATATCTAG GCAGAAACTA CAAGGTTTTA GCCAGTGTG	120
GGCATATCCG TGATTGTAAG AATCCAGTA TGTCCGTCGA TATTGAAAAAT AATTATGAAC	180
GCCAATATAT TAATATCCGA GAAAAAGGCC CTCTTATCAA TGACTTGAAA AAAGAAGCTA	240
AAAAAGCTAA TAAAGTTTTT CTCGCGATG ACCCGGACCG TGAAGGAGAA GCGATTTCTT	300
GCATTTGGC CCATATCTC AACTTGGATG AAAATGATGC CAACCGGTG GTCTTCAATG	360

AAATCACCAA	GGATGCAGTC	AAAAATGCTT	TTAAGAACC	TCGTAAGATC	GATATGGACT	420
TGGTCGATGC	CCAACAAGCT	CGTCGGATCT	TGGATCGCTT	GGTAGGGTAT	TCGATTTCGC	480
CTATTTTCTG	GAAGAAGGT	AAGAAGGGCT	TGTCAGCAGG	TCGCGTTCAG	TCCATTGCC	540
TAAACTCAT	CATTGACCGT	GAATATGAAA	TCAATGCCCT	CCAGCCAGAA	GAATACTGGA	600
CAGTTTGATGC	TGCTTTTAAA	AAGGGAACCA	AACAAATTCA	TGCTTCCTTC	TATGGAGTAG	660
ATGGTAAAAA	GATGAAACTG	ACCAGCAATA	ACGAAGTCAA	GGAAGTCTTG	TCTCGTCTGA	720
CGAGTAAAGA	CTTTTCAGTA	GATCAGGTGG	ATAAGAAAGA	GCSCAAGGCG	AATGCTCCTT	780
TACCCATATAC	CACCTCATCT	ATGCAGATGG	ATGCTGCCAA	TAAATCAAT	TTCCCTACTC	840
GAAAAACCAT	GATGGTTGCC	CAACAGCTCT	ATGAAGGAAT	TAAATATCGGT	TCTGGTGTTT	900
AAGGTTTGAT	TACCTATATG	CGTACCGATT	CGACTCGTAT	CAGTCTGTGA	GCSCAAAATG	960
AGGCGCGCAG	CTTCATTCAG	GATCGTTTTG	GTAGCAAGTA	TTCTAAGCAC	GGTAGCAAGG	1020
TCAAAAACGC	ATCAGGTGCT	CAGGATGCC	ATGAGGCTAT	TCGCTCCGCA	AGTGCTTTTA	1080
ATACACCAGA	AAGCATCGCT	AAGTATCTCG	ACAAGGATCA	GCTTAAGCTA	TATACCCCTTA	1140
TCTGGAATCG	TTTTGTGGCT	AGCCAGATGA	CAGCGGCGGT	TTTTGATACC	ATGGCTGTTA	1200
AATGTCTCA	AAAAGGGGTT	CAATTTGCTG	CCATGGGTAG	TCAGGTTAAG	TTTGATGGTT	1260
ATCTTGCAT	TTATAATGAT	TCTGACAAGA	ATAAGATGTT	ACCGACATG	GTGTGTGGAG	1320
ATGTGTGCAA	ACAGGTCAAT	AGCAAAACAG	AGCAACATTT	CACCCAAACG	CCGCGCCGTT	1380
ATTTCTGAAGC	ATCACTGATT	AAACCTTAG	AGGAAAATGG	GGTGGACGT	CCATCAACCT	1440
ACGCGCCAC	CATTGAAACC	ATTCAGAAAC	GTTATTATGT	TCGCTTGSCA	GCCAAACGTT	1500
TTGAACGAC	AGAGTTGGGA	GAATTTGTCA	ATAAGCTCAT	CGTGAATAT	TTCCAGATA	1560
TCGTAAACGT	GACCTTCACA	GCTGAAATGG	AAGGTAAACT	GGATGATGTC	GAATTTGGAA	1620
AAGAGCAGTG	GCGACGGGTC	ATPGATGCTT	TTTACAAACC	ATTCTCTAAA	GAATTTGCCA	1680
AGGCTGAAGA	AGAAATGGAA	AAATCCAGA	TTAAGGATGA	ACCAGCTGGA	TTTGACTGTG	1740
AAGTGTGTGG	CAGTCCAATG	GTCAATTAAC	TTGGTCGTTT	TGGTAAATTC	TACGCTTGTA	1800
GCAATTTCCC	AGATTGCCGT	CATACCCCAAG	CAATCGTGA	AGAGATGGT	GTGAGTGTG	1860
CAAGCTGAGA	TCAGGGACAA	ATTATTGAGC	GAAAAACCAA	GCCTAATCGC	CTATTCTATG	1920
GTTCGAATCG	CTATCCAGAA	TGTGAATTTA	CCTCTTGGA	CAAGCCTGTT	GCTCGTGACT	1980
GTCCAAAATG	TGGCAACTTC	CTCATGGAGA	AAAAAGTCCG	TGGTGTGGC	AAGCAGGTTG	2040
TTTGTAGCAA	AGGCGACTAC	GAGGAAGAAA	AGATGCTCTT	TTGTCAACTG	TAGTGGGTTG	2100
AAGTCAGCTA	AGCTCGAGAA	AGGACAAATT	TTGTCCTTTC	TTTTTTGATA	TTCAAGACGA	2160

233

TAAAAATCCG TTTTGAAG TTTTCAAAGT TCCGAAAACC AAAGGCATTG CGTTTGATAA	2220
GTTTGATGAG ATTATTGTC GCTTCCAATT TGGCGTTAGA ATAGTGTAAGT TGAAGGGCGT	2280
TGACCAATTT CTCTTTGTC TTTAGAAAGG TTTTAAAGAC AGTCTGAAAA AGAGGATGAA	2340
CCTGCTTAGA ATTGCTCTCA ATGAGTCCGA AAAATTTCTC CGGTTCCTTA TTCTGAAAGT	2400
GAAACAGCAA GAGTTGATAG AGCTGATAGT GATGTTTCAA GTCTTGTGAA TAGCTCAAAA	2460
GCTTGTTTAA AATCTCTTAA TTGGTTAAAT GCATACGAAA AGTAGGGCGA TAAAAATGTT	2520

(2) INFORMATION FOR SEQ ID NO: 15:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 10993 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 15:

TTTCTCGAT AATAACTTCC ACCTATTAT TTGGGATACC CTCTCTTCT TCACCACCAC	60
GTTTCATAGTA GTCATCGGA TAGAGAAAG CTACGATATC AGCGTCTGC TCAATAGACC	120
CAGATTCAG ANTATCAGAC AAGACCGTC TCTGTCCGTG ACCTTGTCT ACACCACGAG	180
AAAGCTGACT CAGAGCGATT ACTGGAACCT TCAATTCTCT GGCTAGTATT TTCAACTGAC	240
GAGAAATTC AGAACTTCT TGTGACGAT TTTCTGACC AGTTCCCGTG ATAAGTTGCA	300
AATAGTCTAT CAAATCAAA CCAAGATTTC CAGTTTCTTG AGCCAATTTA CGAGAACGAG	360
AACGAATCTC TGTAATCCGA ATACCTGGCG TATCATCGAT ATAGATACTG GCCTTAGCTA	420
GATTACCCCTG AGCAATAGTA TATTTTGGCC ACTCCTCATC TGTCATATGC CCTGTACGGA	480
TAGAATGTGA CTCCTAAG CCTTCTGCG CTAACATAG ATCTACCAAG CTTTCCGCAC	540
CCATTTCGAG TGAATAAATA GCAACCGTT TGTCCAATT AGTCCCAATG TTCTGAGCGA	600
TATTCAGGC AAATGCTGTC TTACCAACTG CTGGACGAG TGCTAAGATA ATCACTCTCT	660
CCTCATGAAG TCCTGTGTC ATATGATCCA AATCAGATA ACCTGTGCA ATACCTGTAA	720
TATCGTGTGT TTGTTGCGAG CGAGCTTCCA GATTCCAAA GTTGAGATTC AACACATCTC	780
GAATGTTCTT AAACCGCTT CGATTTCGAT TTTCAGTGAC ATCAATCAAC CCTTTTCTG	840
CCTGAGCAAT AATTTATCA GCTGGTTGTG ACGCTTCTGA AGCTTGGTTG ACAGACTCTG	900
TCAACTTGGC AATTAACGA CGTAGCATTG CTTTTTCTGC AACAACTTA GCATAATACT	960
CCGCATTAGC AGAAGTTGGC ACAGAATTAA CAATCTCAAC CAAGTAAGAC AAGCACCAAA	1020

234	
TATTTCTGTA ATCACCTTGA TTATCAAGGA TAGTACGAAC CGTTGTTGCA TCTATGGCAT	1080
CACCAACGATC GGATAAATCG ACCATGGCTT GGAATAATCAA ACGATGGGCA TACTTAAAAA	1140
AGTCCCGAGA CTCAATGTAT TCTCGCACA AAACAAGTTT ACTCTCATCA ATAAAGATAG	1200
CCCCAAAAC GGATTGCTCA GCTAAGATAT CTTGAGGTTG TACTGGTAAC TCTTCTACTT	1260
CTGGCATCAG ACTTCCCTTC CTTTACAAT CTGTCAAGA AGGTGTAAC TTATCCTCTT	1320
TTACACGAA GATTGATTAC ACTTGTGATA TCTGTATAGA TTTTCACTGG CACATCAATC	1380
AAACCAACCG CTCGAATCG AGCTTGTACT TGAATATGAC GTTTATCAAT CTTAATTCCA	1440
AATTGCTTTT GCAATCTTC TGCAATCTTC TTATGGTAA TAGAACAAA GGTACGACCA	1500
TCTGGACCA CTTTTCACAC AAATTCTACA ACAGTTCTCT CTGCTTCAAG TTGTGCTTTA	1560
ATTGCTTTTC CTCTCGAAT CATCTCAGCG TGAGCTTTT CTTCGGATT TTGTTTACCA	1620
CGAAGTTCAC CTACAGCTTG AGCAGTGGCT TCTTGGCTA GATCTCTTTT GATAAGAAAG	1680
TTTTCGATCAT ACCCTGTTGG TACTTCTCTA ATTTTCGCTT TTTTACCTTT TCCTTTAACA	1740
TCTGTAAAA AGATTACTTT CATCTCTCTT TCTCCTTTTC CTTCATTCA TTTAATACAA	1800
TTTCTGTGAG TTTTTCACCT GCTTCTGACA AGGTATACAT TTTAATTGA GCTGCTGCCA	1860
AATTAAAGTG GCCTCCACCG CCTAATCTTT CCATAATCG TTGTACATTC AGTTTACTAC	1920
GACTTCGAGC TGAGATAGAG ATAAATCCTT GTGTATTCTT CGCAAGAACA AACTTCGCTT	1980
CAATACCTGA CATGCTTAAC ATGGCATCTG CTGCTTACT ANTAACAAT GTATCATAGC	2040
ATTTCAATGTC CTTAGCCTCT GCTATTAGTA CATCTGAACC TAATTTACGC CCTGTAAAA	2100
TAAGTTCATT GACCTCACGA TATTCTTCAA AATCTGTGCG AGCGATTTC TGGATAGCAA	2160
TACTATCACT TCCGCGCGTT CTGAGATAGC TAGCAACATC AAGTGTCCGA CTAGTTACTC	2220
GCGAGGTGAA ATTTTTAGTA TCCAACATCA TACCAGCAT CAAGACACTT GCTTGCATAC	2280
GACTCAAACG ATTTTCTTA GAATTCGGA ACTGAATCAA TTCCGTACC AACTCACTGG	2340
CACTACTTGC ACCACTTTCG ATATAAGTAA TAACCGATT ATCTGGAARA TCCTGATCCC	2400
TTCTATGOTG GTCAATAACA ATGGTTTGGG TAATAAATC ATAAATTTCT TTGTGATAATG	2460
TTAAGCTGTT CTTTGAATGG TCTACAAGAA TCAACAAGA ACGATTGGTC ACCATCCCCA	2520
TGCAATCATT AACAGACAAC AACTTCGTAA CTCCTCTTTT TTCTATGAAT GAAACAGCTC	2580
GTTCATATC TGGAGACATT TGTCTTTCAT CATAAAGAGC ATAGCTATT TTCAATCACAT	2640
TGCTGGCGAA CAACTGCATA CCTACAGCAG AGCCCAAAGC ATCCATGTCT AAATTTTGT	2700
GACCGACTAC AAAAACCTGA TCTACACTCC GAATCTTATC TGAATAGCT GTCATCATAG	2760
CGCGGTACG AGTCCGTGTA CGCTTGATTG AAGCAGCAGA CCCACCACCA AAATAAATCG	2820



GATTTTCGT	TTGTCGTTT	TCCTTAACAA	CCACCTGTC	GCCACCACGT	ACTTCAGCCA	2880
AGTTCAAAT	GAGCAAGCA	ACTTTCCTA	TCTCATCATG	ATTTCCATCG	CCATAAGAAA	2940
ATCCCACT	TAAGGTCAAG	GGCAACTGTC	TCTGTTCGA	CTCTTCTCTG	AAAGCATCAA	3000
TACAGAAAA	TTTATCATTC	ATCAAGCCCT	CAAGCACCCT	GTAGTCAGTA	AATAGATAAA	3060
ATCGATCCAT	ACTTACCCGA	CGAGAAAACA	TCATGTGTTT	TTCTGAAAAA	TCTGATATAA	3120
AATTAGCTAC	AAAACATTTG	ATTTGACTAA	TATCTGACTC	AGAAGTTTCA	TCCTCCAAAT	3180
CATCATAAAT	ATCCACAGAG	ACAATCCCAA	TCACCTGGTCT	ACTTGTATACC	AATTCATCTG	3240
TTATGGCTTG	TTCCCTGGAT	ACATCTACAA	AATACAAAAC	ACCGGAAGAA	GCATCCATAT	3300
GAACAGCATA	ACGCTTCTCA	CCAAGCTTGG	CATAAGTAGA	CGGATTTTCT	ACTGAAGCCT	3360
TGATAATCGT	TTGAACAGCT	TCTAAATCAA	AATCAACCATC	TTCTTGGTTC	AAAATCAATT	3420
CAGCATAGGG	ATTAAACCCAC	TCAACCTCTC	CAGAAGATAA	ATTCAATTTT	ATAACACCTA	3480
CAGGCATCTG	TTCCAAATAGA	GCTGTCAAAC	TTTCTTCCGC	TTGGTGGTPT	ACATACGTGA	3540
TCTGTCTCTAC	ATCACTCCTT	GTATAATGCA	CTCTCAGTTT	CTTAAATAAA	AAAACATAGC	3600
CTCTTACAAA	AAGAAACAAA	ATTAACCAAG	TCAACAGATT	ATTATTAACA	AAAATAATGA	3660
AAGTGGATAA	GACTCCAAAC	GCAATCAATC	CTACTAGAAT	AGGAAAAAAT	GGACTTACAT	3720
AAAAATTTT	CATTCAAAAC	CTCTTGGCAC	CCATATATAC	ATAATACCCC	TCAAAAAGCG	3780
ACTTTTTAAA	AGTGTAATCA	GTAAATCTAT	CAATTATAAG	AAAAAGGTAG	TTTACAATTC	3840
AGTAAACCTA	CCTTACACA	TATTGAATTT	AAGATTCTTT	AACCTCTAAC	AAACCAATTT	3900
CGCATCTCTC	ACGACGATAA	ATCACATTTG	TTGTCTGATC	TTCAACATCC	ACATAGATAA	3960
AGAAATCATG	CCCCAATAAA	TCCATTTGTA	GAATTGCTTC	TTCCAAATCC	ATTGTTTTTA	4020
AATCAATTTG	TTTTGAAGCA	ACAACTTTAG	ACTGGACAAT	ATTTGAATCT	TCCACCAAG	4080
CATCTGTAAA	TAAATTGACCA	GTTGCTACCT	TATTTTTATT	TTTACGCTCG	ATTTTTGTTT	4140
TATTTTTACG	AATCTGACGT	TCAAATTTAT	CAGTTACAAG	GTCAATTTGA	CCATACATAT	4200
CTTGAGATAC	ATCTTCTGCG	CGGAGAGTAA	TAGATCCAAG	CGGAATCGTT	ACTTCCACTT	4260
TAGCCGTTTT	TTACAGATAA	ACTTTTAAAT	TAATTCGGGC	ATCCAATCTT	TGTTCTGGTT	4320
GGAGTACTTT	TTGATCTTTT	TCGAGTTTAG	AAACTACATA	ATCACGAATT	GCTTCTGTGA	4380
CTTCTAGGTT	TTCAACACGG	ATACATATAT	TAATCATATG	AGTACCTTCT	TTCTAAACAT	4440
TTTTTTTTTT	ATGATTTTTT	TATAACGCCT	TCATTCATAT	TTTTGCAAAAT	TTTTCTCTCA	4500
CTTACAAGGG	AAAATGTTTT	TACATCTCTA	GCACCGCTTT	CTTCCAAACG	TTTCTTAAAC	4560

CGATTTATAG TTGCTCCTGT AGTATAGATA TCATCTATAA GTAGGATTTT TTTAGGAATA	4620
GTGACTCCAC TTTTAATAAA GAAAGGAAAT TCTGTCCUCA AGCGCTCTGA ACGATTTTFA	4680
GAAGAACTGG CTCCTCTCTC TCTTTTCTCT AATAAATCCA GATACCTAAA GCCTGCTGCC	4740
TCTACCAAGC CCTCAACCTG ATTAATCTCT CTATTAGCCT ATCTATCAGG ACTTAGGGGA	4800
ATTACAACAA ATTGATACT TTTGACTTT TCGAATCTCT CACTTAAAA TGAAGCGAAA	4860
ACTTTTCTTA ACAGGAAGTC TCCATCAAC TTATACCGAC TGAATAAATC CTTCATAGCT	4920
TGATTGTAA TAAAAATCGC TCTATGACTG ACTTCAACTC CCTCTTTACA CCAAAGTTGA	4980
CAATCTTGAC ACTTTGTTGA CAACCTCTGT TTCATACAAT TTGGACAGTT CTCCTCCCA	5040
ATTCTTTCAA AAGTAGAATC ACAGTCTGAA CAAAGACAAG AGTCATCATT CUTCAGAAGT	5100
AAGAGACTAC TAAAAATTAA AACAGTCTTC ATAGTCTGCC CACATAACAA GCACCTCATA	5160
GACCAGCTTC CTATPTCATC ATCTGAATTT CCTTAATCGC CTCTCTGATT GAAGCATTTA	5220
ACCCATCATG GAAGAAAAGC AAATCTCCGT TCGGCTCATC CATGCTTGGT CCAACTCGTC	5280
CACCAATCTG AATCAAACTA GACTTGGTAA ACAACGATG ATTGGCTCT ACTACGAAAA	5340
CATCCACACA AGGGAAGGTA ACTCCGCGCT CCAAGATTGT CGTACTGATA AGTATTGTCA	5400
GTTCTCCATC TCGAAAAGCT TGTACTTGCT CTAAATGATC CTCTGTACAA GAAGATACAA	5460
AGCCAAATTT CTCAATTGGA AATTGCTCTT GTAAGATTTC TGCTAACTGC TCCCTTTCT	5520
TAATTTCTGA AGCAAAAATG AGTAAACGAT AAGCTGTCTT TCTCTGCTTC TCAATATAGG	5580
ACTTTAACTT TGGTGACAAA CGATTCTTGT CTAAGTAGCG ATTAAAATCC GATAACCAA	5640
TTGGTTTGG AATAATCAAC GGATTTCAT GAAACGCTCT CGGTAAATTC AGTCTTTTFA	5700
GTCTCTCAA ACGGACCTTT TTATCTAACT CATGGTGA AGTGGCTGT AAAAAGATTC	5760
TCATTCGATF CTCTTTTACA CTATTCTTGA CAGCGTGGTA AAGCATGGGA TTATCAACAT	5820
AAGGAAAAGC ATCTACTTCA TCCACTATCA GCAAAATCAA AGCTTGATAA AACTTCAATA	5880
ACTGATGGGT TGTGCAACA ACTAGTGGTG TTCGAAATA AGGTTCOGAT TCTCCATGTA	5940
GCAAGCTAT CCGCAAGAA AATCTCTGTT CGAGGCGCTT GTACAGCTCC AAACAAACAT	6000
CTATGCGAGG ACTAGCCAAA CACACTGCAC CACCCGCAAT GATCACTTTA GCCACTACTT	6060
GATAAATCAT TTCTGTCTTT CCAGCTCTTG TTACCGCATG AACTAAGGTT GGCCTTTGCT	6120
TGTCTACTAC TTGAAGCAAT CCTCTGACA CCTTCTCTTG AAAAGGAGTT AATTGGCGCG	6180
GCCATTGAG AACATCTTGC TTGGAAAAAT CCTCTGCGG AAAATAGTAT AAAGTTTGAT	6240
CACTCTGTAC TCGCTTCATC AGCAAGCACT CTCGACAATA GTAAGCACCG ATGGGCAAT	6300
ACCATTTCTC TAGAATAGTA CTATTACAGC GTTGACAGAA AAGTTTCCC TTCTCTTTC	6360

TCATTGCTGG AAGTTCTCC GCCAACTGAC GTTCTTCTTC TGTTAATCA TTCTCAGTAA	6420
ATAAACGACC GAGATAATCT AAATTACTT TCATACCTT TTAATCGTAA AAACAGCAC	6480
TTTAGATGAT TTTTAGTAC AATTAAATCA TGGAAATTAG GACAAATAAA GAGGACGGTC	6540
AAGTCCAGA AGAAATCAAA AAATCTCGCT TTATCTGCCA TGCCAAAGCT GTTTATAGCG	6600
AAGAAGAGGC TCGTGACTTC ATTACTGCCA TCAGAAAAGA ACACTACAAA GCGACACATA	6660
ACTGCTCTGC CTTCAATTAT GGAGAACGTA GTGAATTAA ACGTACAAGT GATGATGGTG	6720
AGCTAGTGG TACTGCTGGT GTTCCCATGC TTGGGGTACT AGAAATCAC AATCTCACCA	6780
ATGTCGTGT GGTCTGACA CGCTACTTTG GTGGTATTA ACTAGGCGCT GGAGACTAA	6840
TTGCTGCTTA CGCCGGCAGT GTCGCCTTAG CTGTCAAAGA AATGGTATT ATTGAAATAA	6900
AGAACAGGC TGGCATGCT ATTCAATGT CTTATGCTCA GTACCAAGAG TACAGTAACT	6960
TCCTTAAAGA ACATGGTCTC ATGGAGCTGG ATACAACTT TACAGATCAA GTCGATAGCA	7020
TGATTTATGT TGATAAAGAA GAAAAAGAAA CTATTAAAGC TGCACTTGTG GAGTTTTTTA	7080
ATGGAAAAGT CACTTAACT GACCAAGGTT TACGAGAGGT TGAAGTTCCT GTAAACTTAG	7140
TGTAACAAT GAATATPACA CGGTTTCGTT GACATCTCA CAATYACTTT AGCGAGCAA	7200
ATAAAAAGAG GCGTACCAA ATATACTAGA AATGAAGCA ATTCAAAAGA AACCTGATAT	7260
CGTTTTCTT CACACCTATT TACTAGAAAT AGCTGAAGCG AATCACTTGA AATTAATGA	7320
CTTTGATCTA TGATATATAG AAATGGTATG GATAGCGTTA TACTAAAGAT ATCTTATACA	7380
AAGAGGTATT CATATGCTA TTTATAACAA CATTACTGAA TTAATCGGT AAACACCGAT	7440
TGTTAAACTT AACAAATCG TGCCAGAAGG TGCTGCAGAC GTCTATATAA AGCTTGAAGC	7500
ATTTAATCCT GGTTCATCTG TAAAAGACCG TATTGCCCTT AGCATGATTG AAAAAGCTGA	7560
ACAAGATGGT ATTCGAAAC CTGGTCTTAC TATTGTTGAA GCAACAAGTG GAACACCGG	7620
TATTGGACTT TCATGGGTAG GTGCTGCTAA AGGGTATAAA GTCGTCAATG TTATGCCCTGA	7680
AACTATGAGT GTAGAAGCAG GTAAAAATTAT CCAAGCTTAT GGTGCTGAAC TCGTCTAAC	7740
TCCTGGTAGC GAGGGAATGA AAGTGCTAT TGCTAAGGCT CAAGAAATCG CTGCTGAACG	7800
TGATGGTTTC CTTCTCTTC AATTGACAA TCCAGCTAAT CCAGAAGTAC ACGAAAGAAC	7860
AACAGAGCT GAGATACTAG CTGCTTTCGG TAAAGATGGA TTAGATGCTT TTGTTGCTGG	7920
AGTAGGTACT GGTGGAACGA TTTCTGGTGT TTCTCATGCA CTCAAATCAG AAAATCTTAA	7980
CATTCAGGT TTTGCGATG AAGCAGATGA ATCTGCTATT CTATCTGGTG AAAAACCTGG	8040
TCCTCACAAA ATTCAGGTA TCTCAGCTGG ATTTATTCCT GATACACTTG ATACTAAAGC	8100

238			
CTATGATGGT ATCGTTCGPG TAACATCAGA TGACGGCTCTT	GCATCCGAC GTGAATTCG	8160	
TGGAAGAGAA GGCCTCCCTG TAGGGATTTC CTCAGCTCAG	GAGCCATCGA	8220	
GGTTGCCAAA AAATTAGGTA CAGGTAAAA AGTCTCTGCC	CTAGCACCAG ATAACGGTGA	8280	
ACGTTATCTC TCTACAGCAC TTTATGAATT GTAACCGTCC	ATAAGGAAG TCTATTGAAA	8340	
AATCTCCAGA CTAGAGAACT CACGATAGT TCTAATCTG	GAGATTTCTT ATTTGCACCT	8400	
TTCTTGACA ACTTTAGTCC ATGGTAATA GGCCTCTAAA	ACCTCTTGT TTACGAGCT	8460	
TTCCAOGTTT GGAAGACATT CTAGAAGATA GGATAGATAT	TTCTCACTAT TTATAATGGA	8520	
TTGAATAAG ATATGAACAA ATCGATTAGA ACATGATGGT	AAAGCGTAAT CCCTTGTTTC	8580	
TCAGCTTTCC CAGACAAAA AGTCCAATAG TAAGTCAGTC	GACTATCACT CTCTAGCAAC	8640	
CTATAAGAG TTTTCATCCGC ATGAAGTAAG GGCTGAGTCA	ATAGTCTCTC TCGCAAGAGG	8700	
TTATAAAGGG GCTCCAATA GTATTGACTC GTCTTGATAT	GCCAATTAGA GATTTCCTTA	8760	
CGTGTGATTG GTAAACCCAT CCTGAGCCAA TCTTCTCTT	GGCGATAATT GGGTACCTTC	8820	
AGATTAAACT TCTGATGGAT GGTGTGAGCG ATAAATGAAG	CTGAGCCAAA GTTATGCGCT	8880	
AAAGGGGCTT TAGGAATAGG AGCTTTCACA AGCTTATCCA	GATGATTATC TTTTACTCGT	8940	
TATGGCAAT GCTATATGGC ATAAATCAAG TACCTTAAG	ATTCGACTA ATATTGGCTT	9000	
TGCATTTIAT CTCCATACA CACCAGAGAT GAACCCCAT	GAACAAGTGT GGAAGAGAT	9060	
TCGTAAACGT GGATTTAAGA ATAAAGCCTT TCGAAGTGT	GAAGATGCA TACAGGACT	9120	
GGAGAAGGAG GTGATAAAGT CCATCGTTAA TCGGAGACG	ACTAGAATGC TTTTGAAGA	9180	
CAGATGAGTA TAAAAAGAAA GTCTCTATT CAATGAAAT	CACGACTTTC TGATGAATTC	9240	
ATATGAAAT GAAATAGAAA CAGGATAGTC AAATCGATTT	CTAACAAATG TTTAGAAGCA	9300	
GAGGTGACT ATTCTAGTTT AAATCCACTA TATTTGGGA	GTGATAGAAA AGCCCTTCAT	9360	
CAGCCAATCT ACTGTGTCAG GTGCGAGAGC TTTGACATCC	TTTTCTGTAC TGGACCAAGT	9420	
CAGTTTCCG TTCTCAAAGC GTTTATATAA TATCCAAAT	CCTTGACCAT CCCAGTAAG	9480	
AACCTTAAAG CGGTCTTTAC GTCCACCACA AAAGAGAAAG	ACTTGATCGG AGAAAGGATC	9540	
CAATTCAAAG TGGGTTTAA CTACATAGGC TAATGAGTCT	ATTCCTGCCC TCATATCTGT	9600	
CTTGCCACAA ACAAGGTGAA CTTGACCTAA ATCACTTAGT	TGAATATCA TAGTACAATA	9660	
CCTTTCTCC CATAATATT TTTTATCTGG TATACTGAAA	GTGGGGGAAT TAGGATAGAT	9720	
ACCTGTGTAT GACGCGCTTA CTATGAATTT GAAGTATAGT	CTCCTAAATG CACTTAGCCC	9780	
TTATTATAGG GCTTTTGTGT TTAATTTATTC TAATCGAGTG	AGACTGGGGA AAAACAAT	9840	
TCAGAAAAA TCTAAGCCCT ATACAAAAA GGAAGCAAT	TGCTTCCTTT CTATTATTAG	9900	

239

TTATTCAAGG CTGCTGCUAT TGTAGCTGCA ACTTCAGCTT CGAAGTCGTT TGCAGCTTTC 9960  
 TCGATACCTT CACCAACTTC AAAGCGAGCA AACTCAACTA CCGAAGCGTT AACTGATTTCA 10020  
 AGGTATGCTT CAACCTGCTT GCTGTCAATCC ATGATGTAAA CTGTGCAAG AAGTGTGTAA 10080  
 GCTTGGTCAA CTTTAGTGTT ATCAAGCATG AAGCGATCCA TTTTACCTGG AATAATTTTG 10140  
 TCCCAGATT TTTCTGGTTT GCCTCTGCA GCCAATTCAG CTTTGTGTC AGCTTCAGCT 10200  
 TGAGCAATAA CATCATCAGT TAAITGAGCT TTTGATCCAT ACTTCAAGTG TCGAAGAGCT 10260  
 GGTITATTTAA CCATTGACG GCTTTCGTTG TCTTGGTCGA TAACGTGATT CAATTGTGCC 10320  
 AACTCATCTT TAACGAATTG CTCATCCAAT TCTTGTAAAG AAAGAAGCTT TGGTTTCATC 10380  
 GCTGCGATGT GCATTGACAA TTGTTTAGCA AGTGCTTCGT CTCACCTTC AACAACGTAA 10440  
 ATAACACCGA TACGTCCACC GTTATGTTGG TATGCTCAA AGTGTGTGC GTCTGTTTTT 10500  
 TCAATCAATG CAAAGCGAG GAAATGAGATT TTCTCTCCGA TAGTGTGCT TGCAGATAG 10560  
 TATGCAGCTT CAAGAGTTTC ACCTGAAGGC ATTATCAAAG CAAGAGCTTC TTCGTTGTTA 10620  
 GCAGGTTTTT CTTCAGCAAT GACTTTAGCT GTAGTATTTA CCAATTCAC GAATTGAGCG 10680  
 TTTTTTGCAA CGAAGTCAGT TTCACGTTT ACTTCAATAA CTGCTGCAAC ATTACCGTTA 10740  
 ACATAAACAC CAGTCAACC TTCTGCAGCA ACACGGTCAG CTTTCTTAGC TGCTTTAGCC 10800  
 ATACCTTTTT CACGAAGCAA TTCAATCGCT TTTTCGATGT CACCGTCTGT TTCTACAAGC 10860  
 GCTTTTTTAG CGTCCATAAC ACCGGCACCA GATTTTTCAC GCAACTCTTT TACAAGTTTA 10920  
 GCTGTAATTT CTGCCATTTT AATTCTCCTA TATTTTTTGA AAATAGGAGA GCGCGCTAA 10980  
 GCCCGCCTC CGG 10993

## (2) INFORMATION FOR SEQ ID NO: 16:

## (i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 841 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 16:

CGACGGGAG GTTTGGACC TCGATGTCGG CTCGTCCGAT CCTGGGGCTG TAGTCGGTCC 60  
 CAAGGGTTGG GCTGTTCCCC CATTAAAGCG GCACGCGAG TGGGTTCAGA ACCTCGTGAG 120  
 ACAGITCCGGT CCCTATCCGT CGCGGGCGTA GGAATTTGA GAGGATCTGC TCTAGTACG 180  
 AGAGSACCAG AGTGGACTTA CCGCTGGTGT ACCAGTTGTC TTGCCAAAG CATCGCTGGG 240

		240	
TAGCTATGTA	GGGAAGGAT	AAACCTGAA	AGCATCTAAG
GAGATTTCCT	ATGATATAT	ATCAGTAAGA	GCCCTGAGAG
GAGTGGGAAG	TGTGGCGACA	CATGTAGCGG	ACTAATACTA
AAGTAACATGA	GAATATGAAA	GCGAACGGTT	TTCTTAAATT
GTAGGTATTA	CTCAGAGTTA	AGTGACGATA	GCCTAGGAGA
ACACAGAAAT	TAAGCCCTAG	AACGCCGAA	GTAGTTGGGG
GAGTGCCTT	AGCTTAAATC	CGCCATAGCT	CAGTTGGTAG
TGATGTCGTA	GGTTCGAGTC	CTACTGGCGG	AGTAATCGAT
TCCTCTTTTT	GTATCAATTT	GTATCACCAG	GCATTTTCAT
GAGAACTTTC	TTTTTTTCCA	TGTGCAATCC	AAGTTTGGCA
TTAGATAGAT	GCTACTATAT	TCTAATTCAG	TGGTATTTAG
TTTGCAATTC	TGTACTAGC	ATGATATGAA	GTTTATTTTC
TAGTCCCAT	TTGAGAAAGA	AGGGCAGCCA	GAAATGGTTC
CTTCTAAAT	AGCGTCTCTT	TTGTGATGAG	CATGTTTTTG
GGAAAGCTT	GCTTTGATAG	TGCTCAATCA	TATCATACTT
AGCTGGAAAG	ACTAATTCCT	GCTTTTCTTA	CTAATTTGAC
GCTGTCCAT	CAGTAATGT	ACCATAGCAT	TTCAATAGT
TGTTACTTTC	TGTCATATTT	CCTCTTGTA	AACAAATTAG
TTTTTTATCT	TGTAATTTAG	ATTTTAAAT	GTATACTCTA
ATATGTTTAA	AAAGGAGAA	ACTAAGTTTA	AAGAATGAA
CCTTTATTAT	TGTCATGATC	GGGATTTCTC	TTATTTCCAG
TGTCATCAAT	GTGGGATCCA	TATGGGCAAT	TGTCTGACTT
ATGATAAAGA	GGCTTCCTAT	AATGGTAATA	CTATGGCAAT
ATTTAAAGA	AAATAAAACC	TTGGATTTTC	ATTTTGTAGA
GNTTGAAGA	TGGCGATTAC	TATATGGTAG	TGACTTTACC
CAACTACATT	ATCCAATATT	CAATCGACAG	CAGCTTATCA
AAACTGAGAT	AAGTATTTCT	GTATCTCAA	ATTCACCTGA
CAATTGTAGC	TTTAGTACAA	GATTTACAGG	GAAGTTTAGA
CTAATCTTTC	GACTTTAAAA	AATCAATCTA	ATCAAGTATC
TGATAGGATT	GTCAAGTGGG	TTAACAGAGA	TACAAGGAGA

300

360

420

480

540

600

660

720

780

840

900

960

1020

1080

1140

1200

1260

1320

1380

1440

1500

1560

1620

1680

1740

1800

1860

1920

1980

2040

CTGCCAGTCA	GTCGATTGCA	TCAAGTGTAA	ACGCATATAC	TACAGTGTGT	GATAAAGTPT	2100
CTCAGGGCGC	AAGTCAACTA	AGTGAAAAAA	ATGCCACCTT	GACAGGTAGT	TTGGATAAAC	2160
TAGTTTCAGG	CTCAACACCC	TTGACACAAA	AATCTTCTAG	ATTGACAGCA	GGAGTTGGTT	2220
AATTACAAATC	AGGATCTGGG	CAATTAGCAG	ACAAATCCAG	TCAGTTACTT	TCAGGTGCTT	2280
TCCCATTAGA	GAATAGAGCT	AATAAATTGG	CAGATGGATC	TGGGAAACTA	GCAGAAGGTG	2340
GAACAAAGTT	AATCTTGTGA	TTGGAAGATT	TACAGACAGG	ACTTGCTTCT	TTAGACAAAG	2400
GACTAGGTAA	TGCTAGTGAT	CAACTCAAAAT	CAGTATCAAC	AGAATCTAAA	AATGCAGAGA	2460
TTTTGTCAAA	TCCACTCAAT	CTTTCAAAAA	CAGACAATGA	TCAAGTTCCT	GTAATGGAA	2520
TGCAATAGC	TCCTTATATG	ATATCAGTTG	CTCTTTTTTT	GCAGCAATAT	CAACAAATAT	2580
GATATTTGGG	AAATTGCCTT	CAGGACGTCA	TCCAGAGAGC	CGTTGGGCTT	GTTTGAATC	2640
TTGAGCTGAA	ATAAATGGTA	TTATAGCTGT	TTTGGCAGGA	ATTTTGGTAT	ATGGAGGAGT	2700
TCAGCTTATT	GGTTTAACTG	CTAATCATGA	GATGAGAATA	TTTATTCTCA	TCATCCTAAC	2760
AAGTTTAGTA	TTTATGTCTA	TGGTGACCAC	TTTAGCAAGC	TGGAATAGCC	GTATAGGAGC	2820
TTTTTTCTCA	CTTATTTTGC	TTTTACTACA	GTTAGCATCA	AGTGCAGGTA	CTTATCCACT	2880
TGCTTTTGACA	AATGATTTCT	TTAGATCTAT	TAAATCCCTGG	TTACCAATGA	GCTATTCACT	2940
TTCCGGGATTA	CGACAAACAA	TCTCTATCAA	CAAGTCATTT	TCCTAGCTGT	CATACTAGTT	3000
CTATTACTA	GTTTAGGTAT	GCTAGCCTAT	CAACATAAGA	AAATGGAAGA	AGATTAAAAA	3060
AATCGACCGA	TTAAGTGGTC	GATTTTTTAT	GCCTTAGATG	ACTTTCGTCT	GTGATTATAG	3120
ATTCCAAATA	GTAAGAGAGA	AGTAAAGGAA	CAGATTGCTC	CAGTAATATA	ACCATTGGGA	3180
ATGAAGGAAA	GTGTAATAGT	TCCTTTCCCC	TTGGGAATGT	CAACTTTTCA	AAATCCAGTT	3240
TGAGCTTGTT	TAAATTTCTAT	TTTCTTACCA	TCTTGGTAGG	CAGACCAACC	TTTGTCTATA	3300
GGAATGGTGA	AGAAAATAGA	TGTATCTTGT	TGGACATCAT	ATGTAGCAAA	AACCTTGTTT	3360
TTAGAAGITG	ATACTGTGAC	AGGTGTGTTCT	TTAATTTTTT	GAAATGCCTC	GGTGAAGTT	3420
TTGGTATCTA	AACGATAGAA	GGTAGGAGAT	TCAAATGATA	CTTGTGAATT	TCCAGGGAAA	3480
CTAACATTGA	TATTGAAAGT	TTTTTTCTCT	TTAGTATATC	CTAGATTAAA	GAGGAGAGAG	3540
ACATTATCAG	TTGTAAAAGT	CTTTTTTTCA	CCATTACAAA	GGATGTCAC	CTTCTTTTGT	3600
TTATCGTTAG	AAAAGTGAGG	GTTTATGAAA	GAGAGATAAA	CTTGGCTGTT	TTCTGGAAGT	3660
TCAATTTGAT	ACTGGATTGC	TGCATCTTCA	TTTGAAGAAC	TTGTGACACT	AATCAAAATCA	3720
TTAGTATTTT	CTATTTTTTC	TGTTTTTTTCA	TAAGGTATTG	GAGAAAAATA	ATCAAAATTG	3780

ACGTAGCAAA GTTGATTAA AAATGAGGUC TGATTATCCA AGGTATGTC ATPGAACCTG	3840
ACATCATTTGT AAACAGATTG ACTCGCAACT GCAATCGAAA GAGAGTATPG ATTTTCATAT	3900
AGGTAAGAT TATCTTTTTG ATAGATATCT TTAAGCCAT ACTTATCAAT AGGACTGTCT	3960
GAGATATTGT ACTGGATACC AAATAAACTA TCAGCCAAAA TACTATTATT TGCATATCGG	4020
AGATTGAGAT TAGTCCCGAA GGATTTAAAA CCAAGTTTAT CTAAAGTAGA GCTTGATGAA	4080
CGATTTCGAA CAGATGAAAA TTGAGAGATT CCAATTGAGT TGAATTTCAT ACTGTCAATTT	4140
CCTGTCTGAG TTTGAGTTT TTAGTACGA GTAAATTGAT TTCCAATATA TGTGAGAAAA	4200
GATTCCATAG CTGGGATATC TCGACTATAA GCACCTCGAG AAGCAATCC CCAATCCCTTA	4260
GCAATTCCGT CCAATTGAGA TGAAGCATT AAACCTATT CAACGAGAT AAATAAGAG	4320
ATTAGAATGG CAAATAGATT CACAGATATA AACTTTTTGA TAACCTCAAG GAGTAAAAGA	4380
GAATAGACAA CCAAAATTC AAGAGTAAGC AGAATATTCA AATCTGTAA AAAAGAATAA	4440
TGCGATTTTA GATAGATGTT AGCTAAAAAT CCTGCTACTA CAAGAAAAAG CGAAACTAAA	4500
AAATTCGAGA CTTTAAGTTC TTTCAGACGC TTTAAGACTT CTGCTGCTGT GTAAATTAACT	4560
AAGGTAGAGA AAATCCAAGC ATAGCGATGT AAAACATGT TTGGAGTATG CATGCCCTGC	4620
CAAAATTAAGT CAAGAGCTTC TATGTAAAAG CTGCAATTA GAAATGCCAA GAATATTACA	4680
TATATGAGTT TCACGTGAAA CTTAATAGAT TTACGCGTAA AAAATAAAT GGTCAAAATA	4740
AAGGAAATTA GTCCACAAA AATCATTTGG ATGCCCCAT ACTTTGTTGT GTCAAAGGAA	4800
CCAATGAATT CTTTAGCAAA GAGATCAAGA TACCAGCTAC TTTCACTTTG AAACCTTTGA	4860
ACTTCAGTCA ATTTTTCCCC ATGTGCTGT AAATCAATA GAGTGGGAAG AGTCATTAATC	4920
AAACTAGCCA TACCAGCTAA AAAGGAGATA ACTATGAAAT CAAGAACAGA TGATTTTCGA	4980
GTCTTAAAGT CCCACGAAAT TTGACAGAGA TACCAGAAAA TAAGAAACAA TACTGTCATA	5040
TATCCAAAAT AATAATTTTG AATAAATAAG ATGACAGAC TTGTAAAGTA CAATAGGAGT	5100
TTCTTTTCAG TTAATCAGT ATGTAAACCA GTTATAATTA AAGGAATCAA GATAAAAAACA	5160
TCTAGCCAGG TTTTATCTC TAATTGACTG ACAGTGAAC TCATCAGAGC ATAGGAAGTA	5220
GATAAGGCTA GTTTTAAAT CTGAGGGATA GATTGAACAA ATTTATTCAA ACTAAAAAG	5280
GTTGACAGAC CAATCAATCC AAATTTTAAG AGAGTTGTCA GATAGATAGC ATCTGGCATA	5340
TTCTTTAGAT CAAAAAGTA AACCAGAGGC GCGAGAAAC TACCCAAGTA ATAAGTAGAT	5400
AGGGCATAGA AGTTTAGCCC TAGACCACTT GTAAAGGTGT AAAACAGATT ACTATTTCCTA	5460
TGTAGGATAT TTCGTAAAGC TACATCAAAA ATAAGTATT GATGAAGCC ATCTCCYBAT	5520
AGAGGAGAGT TGTGCTATT CCAGTAGATA CTTTGAGATA GATATACTCC AGACATAATC	5580



ACTACAGGAA TGATGAAAGA AATAAAATAG GTTCGATATG TTTTAAAAA TGATTTTCATG	5640
TTACCTCGTA GAATGATAGA AAACCTCAGTT GGTAAACCCA ACTGAGTTTT GAAGTTTAT	5700
TTAGTCTTTC CAAAGTTCTT TAACTTTTGC TTGTACTTCT GCATTTTCTA GGAATTCATC	5760
GTAGTCTTCA TCGATACGGT CAATGACGCC ATTTTATAGT AAGACAATGA TATGTTTAGC	5820
CAAAGTTTGA ATAAATTCGT GGTTCATGGCT GCGAAGATG ATTGATTCTT TAAAGTTTTT	5880
CAATCCATCA TTCAAGCTTG AGATAGATTC CAAGTCCAAG TGATTTGTGT GATCATCAAG	5940
TACAAGGACA TTTGATTTTA AGAGCATGAG TTTTGAAAGC ATGACACGAA CTTTTCTCC	6000
CCCTGACGAG ACATTTCACG GTTGTGTAAC TTCATCTCCA GAGAAGAGCA TACGCCGAG	6060
GAAGCCACCT AGGAAAGTAT TGTCTCTTC TCTTTACTT GCGAATTGAC GCAACCAATC	6120
AAGAATTGAT TCTCCTCTG CAAATCAGC TGAGTTATCT TTTGTAGGT AAGATTGACT	6180
AGTTGTAACT CCCACTTGA CAGTCTCTC ATAGTCATTA TCTCCCATGA TTGCACGAAT	6240
TAATGCAGTC GTTTGAATAT CATTTGTGCC AATAAGTGCT GTCTTATCAT CTGGACGCAA	6300
GATGAAGCTA ATATTATCCA AGATAGTTTC ACCATCAATC TTTACAGTTA AATTTCTAC	6360
TGTCAAGAGA TCATTACCAA TCTCAGTTC CGCTTTAAG TTGATAAATG GATATTTACG	6420
ACTAGATGCC ACAATCTCTT CTAGCTCAAT CTTATCAAGC ATTCTCTTAC GTGATGTTGC	6480
CTCCCTTGAC TTAGAAGCAT TGGCAGAGAA ACGAGCAACA AATCTTTGCA ATTGTTTAAT	6540
TTTTTCTTCT GCTTTAGCAT TACGCTCTGC TAGCAATTTA GCAAGCAAGCT CAGAAGATTC	6600
CTTCCAGAAG TGTAGTTTC CGACATAGAG TTTGATTTTT CCAAGTCAA GTCCGCCAT	6660
GTGATACAAA ACTTTGTTTA AGAAGTGACG GTCTGTGGAT ACTACGATAA CTGTGTTATC	6720
AAAGTCAATC AAGAAGTCTT CPAACCAAGT AATCGATTGG ATATCCAACG CGTTAGTAGG	6780
CTGTGCCAAG AGAAGACAT CTGGTTTACC AAAAGTGCT TTGGCGAGGA GAACCTTTAC	6840
TTTTTACACG TTGGCCAATT CGCTCATGTT TTGTAGTGT AATCTCTCG GAATGTTTAG	6900
GTTTGAAGT AGTTGAGAG CTTCACTCTC TGCTTCCCA CCTCCAAGTT CGGCAAACTC	6960
TCTCTCGAGT TCGGCAGCAC GAACCCCTC CTGCTCTGAG AAATCTTCTT TCATGTAGAT	7020
AGCATCTTTC TCTTTCATGA TGCTATAAAG TTTTTCATTT CCCATGATAA CGACATCAAT	7080
GGCAGCTTCA TCTTCGTAGT CAAAGTGATT TTGACGAAGA ACAGAGAGAC GTTCATCTGG	7140
ACCAGAGAG ATGTGACCAG TAGTAGGTTG GATATCTCCA GCTAAAATTT TTAATAAGGT	7200
TGATTTTCCG GCACCATTAG CACCGATTAA TCCGTAACTA TTTCTTCTG TAAATTTGAT	7260
ATTGACATCA TCAAAAAGTT TCGGATCACT AAAACCTAGT GAAACATCAG ATACTGTAAG	7320

244

CAATGTTTTT CTCCTATATG TGTAAATATAT TTATTCTACT AGAAATACA GAAATATTCA 7380  
 AATTTTTTAT TGTCAATTTT GTGTAAATTA TATTACAGT ATCCTTTTACA CAAATCTGTA 7440  
 AAAAGCAAGG CTGATTTTAT TTGATAAATT ACGGTTATTT CATTAACAAA ATGCTATAAT 7500  
 TGAAGGACT ATATCGAAGG AGAACAAAA GACTAAACCC ATTATTTTAA CAGGAGACCG 7560  
 TCCAACAGGA AATTCGATA TTGGACATTA TGTGGAAAT CTCAAAATC GAGTATTAT 7620  
 ACAGGAAGAG GATAAGTATG ATATGTTTGT GTTCTGGCT GACCAACAAG CCTTGACAGA 7680  
 TCAATGCCAA GATCCTCAAA CCATTGTAGA GTCTATCGA AATGTGGCTT TGGATTATCT 7740  
 TGCAGTTGGA TTGATCCAA ATAAGTCAAC TATTTTATTT CAAAGCCAGA TTCAGAGTT 7800  
 GGCTGAGTTG TCTATGTATT ATATGAATCT AGTTTCGTTA GCACGTTTGG AGCGAAATCC 7860  
 AACAGTCAAG ACAGAGATT CTGAGAAAGG ATTTGGAGAA AGCATTCGA CAGGATCTCT 7920  
 GGTCTATCCA ATCGCTCAAG CAGCTGATAT CACAGCTTTC AAGGCTAATT ATGTTCTGT 7980  
 TGGACACAGT CAGAAACCAA TGAATGAGCA AACTCGTGA ATTGTTGTT CTTTTAACAA 8040  
 TGCATATAAC TGTGATCTCT TCGTAGAGCC GGAAGTATT TATCCAGAAA ATGAGAGAGC 8100  
 AGGGCGTTG CCTGTTTAG ATGGAATGC TAAATGTCT AAATCACTAA ATAATGTTAT 8160  
 TTATTTAGCT GATGATCGG ATACTTTGCG TAAAAAGTA ATGAGTATGT ATCAGATCC 8220  
 AGATCATATC CGCGTTGAGG ATCCAGGTAA GATTGAGGA AATATGTTT TCCATTATCT 8280  
 AGATGTTTTT GGTGCTCAG AAGATGCTCA AGAAATGCT GATATGAAG AACGTTATCA 8340  
 ACGAGGTGGT CTTGTTGATG TGAAGACCAA GCGTTATCTA CTTGAAATAT TAGAACGTGA 8400  
 ACTGGGTCG G 8411

## (2) INFORMATION FOR SEQ ID NO: 17:

## (i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 9064 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 17:

TGCCGTACTC AAGTACAGCC TGGCTAAGT TTCTAGTTT GCTCTTTGAT TTTCATTGAG 60  
 TATTAGTAAC CAAATACGA CCACATAGCC AGCCCCATG AATATAGCCA TTAAGCTAG 120  
 CATGGAATTT AGGAATTA AAACCAACGC AGATACAAAG GTTAGCACA AAACATTAAA 180  
 AGCAATGGTG TCAGAGCCA AGACTAGAT ATAGGGGTGC AACCGATCTA AAGTTTGGGA 240  
 ATCTAGGAAA AATAAGTGT TATACATGAT GACCTCCTCT ATGGCTGAAA AGCAAGCCTT 300

TTGTTTITTT	ACCCCAAGAC	CCTATGTAGA	AAAGTGAGCA	AAAAACGGAA	GSTCGCTACA	360
ATATTTATTGA	TCACATGCAC	CGCATAGGAT	GGATAAATGC	TCTTGGTATA	GCGGGTCAAA	420
CCAGCAAGA	TGATTOCAAC	TGTTGCAAG	ACGAAGATAT	CTAACAGACT	AGGCAGGCTT	480
GAAAAATGAG	GGAGAGCAA	TAAAAAGAA	GGAAGAGCA	AATCAAGACC	AAATCGCGAA	540
TGCTTAAAGA	AAGCATGTGT	CAGTAATCCT	CTATAAATCA	ATTCCTCCAT	CAGTGGAAAC	600
AGAAAGAAAC	GGGCTATATA	AATACCTAGC	TCTGCAAAAT	TAGTCCCACT	ATAACCAATC	660
AATACAGCCC	AACCTTCCGC	AGTTGACTGA	ACATGTTTAG	CTGTCTGAAC	GTTAAAAAGAG	720
ATCTGGAACA	CTAGCACTAA	TACTGTCAAA	ATCGAATACC	AAAGCCATTT	TTTCTTTGGA	780
ATCGGGAAGA	GATAACCATG	GCCTGTCTTA	ACAAGAACCA	CAATCATGAC	TCCAATAAAA	840
AGTAAACTCA	AGATATTTTG	AATCCAGAA	AAATTGCCTA	TCTGAGAAGA	AAATGCCCCA	900
TAGTTTGGGA	CGATAAGCGT	CAGCTGAGAA	AGACTAAATA	CGAAAAATAA	GTAAGAGAAG	960
ACTGCACCTA	TTTGTGAATG	AAGTTGATAC	TTTTTCATAG	AAATCCTCCC	TACTATGACC	1020
TCACCTTGTC	AGGCTCTACT	GCTGTAAGAT	TAAGAAGACA	GTTTGTTTTT	TTTAAAGCTA	1080
ACCTGACTAC	TAGATAATAG	ATACATTAAG	GCATTAAAGA	CAATGAAAAAT	ATGTCCATAG	1140
AATAAAATCA	ACCTCGCATC	CAAAACCAAG	TAAAGTTTGA	TTATCAAAAA	GATGAGCAAA	1200
AGAAATTTGAA	ACCATAAGGT	TTTTTCCAAA	ATAAATTTAA	AGCGATTTCG	AATATCTACT	1260
TCCTTGATTT	TTACCGCCAC	CCCTTTIATTA	GCAAGAAAGGA	AAACTCCTGC	TTCAAAACAAA	1320
CCACTGTAAA	GAACAAGCCA	CCCATAGAT	ACGATAGAGA	TTTGTA AAAA	TGTCCCTPAA	1380
AGAATATCCA	ACACACTACT	CAAGAAAAAT	ACAAAAAATA	ATCTGTATTT	CATATTAAT	1440
ACCTCCATTC	ATTTATTICA	CTAACAAATT	AATAGAGCCT	TCTACTCAAA	TATCCTGTCA	1500
GAAAAGGATA	GAAAGCTACT	TTTTATAATA	CTTCAAGCCC	CACATAGACA	GAAGCCTGAT	1560
AAACAAGCAG	AGAAATACAC	TATATAAGCG	ATTAGTTGTT	GATAGAAATC	TGTTTCTGAA	1620
ATACCTCTAT	ACAAACAAAT	GACAAACATA	AAATCTGCCA	AGCCGATAAA	CMTAAGTTGA	1680
TTGGTCTTAG	GACTAAACAA	ATCATCATTT	ACTTATATTT	AAGAGTATCT	CTTTTATTTT	1740
AATGTATGTT	AGCACTGAAA	AGCAAGACAG	GCCAATAATA	TTTAAATGTA	ACAGTAACGG	1800
GGTTAAGTCT	CTAAAAAAAT	TATCTACTGA	CACACACAAG	AATACTATAC	ATATTTATAGT	1860
CGAAACTATC	TTTTTCTTAT	CCATAATTAT	TTACTCCTTT	CCTAACAAAT	CCAGCTTATC	1920
AATCAAGAGC	GATTTTTTAAC	ATAATGTAGC	AGCACCCGTT	GCAACTTTGA	CAAGTTTAGT	1980
ATATCATTTGT	TTTTTAAAAAT	TTTTTATCCA	AATCTTTGAAT	TGTATCTGAA	ACATCTTGAA	2040

	246	
TTGTTAAAAA ATTTAAAAAG TAAGCATTTAA AACATCTCTT TCCTCTTTAT ATTGTATTGA	2100	
TACCAACTTG TTTGTAGACT TTTTCATCTG CTATCACATA TCATTTTGAC AGGCGAATAA	2160	
ATATTAAAGA AACTCCCCTG TAAATTAAGC TAGCAAATAC AGGGGAGAAA TTTTATTTT	2220	
AGAGAGTACT ATCCGTATCC TTTTGGAG ATTTTGAATA TATTTTCTTA ATTAAGTCAT	2280	
CCATATAAGG ACCAAATATA CCAACTACTA AACCAATAAT AAACTTTTA AAATCCATAA	2340	
TTACCAACAA CATATTGCTG CATAGGCTAC ACCTCCAAGT ATAGCTCCAC CTGCAGCACC	2400	
AGTTACACCT ATTCTATAG CAATGGTCC CAATAGAAAT GTCAAACCGT TGTGCAACAC	2460	
CCATCAATTG CGCCATATGC AACCCCTGCT GCACAACATA TTTTCTTCC CCAATCAATA	2520	
TCTCCACCTT CAACGCAAGC AAGCATTTCA TTATCCATAA CTGCAAAATG TGACATCATT	2580	
TTTGTATCCA TATAGTGTAT CACTTTTCAG TTACGGAACA AGTTAATAT AAAAAATTATC	2640	
AAAAAACAT AGGCAATAAA GAGAAAAATT AATTTATCAT AGATTAGAAA TAATATGACA	2700	
AAACATTTCA ATGATGTTAA TTCAATAGTC TTTTGTTTT TATCGAGAT ACTTATGGAT	2760	
AGATAAATAA GATAGGTTTG AAAAGCGAAG AGAATAATAA AGAATATAGC CTTCATAAAA	2820	
TTTAGCTTTC ATTTTATGA TGTAGCGGTA TAGGCTAAAT ATCCACAAC CACTGCTCCT	2880	
CCAAATTCCTC CTATTGCGC GCCCATGGT CTTAGAAGTC TCCCATATTT CACTCCACCC	2940	
GCTGCACAAC CTAAAGCAGC AACTACAGCT GCTCTCCGG AATTACCTCC ATAAACCTCA	3000	
CTCAGCAATT TTTCAATTTT ATTACAATAA GTATTCTATC AAGTCTCCTT TTAATAAAA	3060	
CCACCCGTTG CCCCTGTTAC TCCTGCCCAA AGATCCACAC CAATTTTAGC TCCTATGTAT	3120	
CCACATGCTC CCATAAATGG TGCTCCAACA CCACTGCGAG CACAAATAGC TGTCCCTAGT	3180	
CCCCAGCCAC CAAAAGCAGC ACCACCACCT TCTAAGACAT TAGTTTGCCA ATTATTTCTTG	3240	
CCTCCTTCAA TACTAGATAA CATAGTTATA TCCATTTCAAT GAAATGTTC CATAATTTT	3300	
GTATCCATGA CAAATCTCT TTTTATTTT TATTTTGTG CTTGTTGTAA CTTTGACAG	3360	
TTTAGTATAT CATCGTTTTT TAAATTTTT CATCCAGATT TTGAATAGTC ATCGAAACGT	3420	
CTTGAATTGC AAAAATTACA TTGACTTCC TGCAAAACTA GAATCCTAGT TCATGATTGA	3480	
TAAATACCAGC ACTCAAATTC ATTCGTAAAT CGAAGCGTTT ACGATGACTT CGATAGGTTG	3540	
TTGAAAACAT TTTAAACGTT TTTACTTTGG CAAAGATGTT CTCAACTTGG CTTCTCTCCT	3600	
TAGATAGCGC ATGGTTACAG GCTTATCTTT CAACCTGTAG CGGTTGAGT TTGCTGGATT	3660	
TACGTGAAGT TTGCTGTTGA GGATATATCT TCATGAGCCC TTGATAACCA CTGTCAGCCA	3720	
AGATTTTACC AGCTTGTCCG ATATTTCTGC GACTCATTTT GAACAACCTC ATATCATGAC	3780	
AATAGTTTAC AGTGATATCC AAAGAAACAA TTCTCCCTTG ACTTGTGACA ATCGCTTGA	3840	

TCTTCATAGC GTGAAATTC TTTTACCAG AATCATTCGC TAAITCTTTT TTTAGGCGCA	3900
TTGATTTTAA CTTCOGTCG ATCAATCAT ACCGTGTCCT CAGAACTGAG AGGAGTTCTT	3960
GAAATCGTAA CACCACTTTG AACAGAGTT ACTTCAACCC ATTGGCTCCG ACGGAGTAAG	4020
TTGCTTTCTG GAACACAAA ATCAGCGCA ATTTCTTCAT AAGTCCGGTA TTTCTGCACA	4080
TATTGAAGAG TGGCCATAAG AAGTCTTCT AGGCTTAATT TAGGTTTTCG TCCACCTTTT	4140
GCSTGTTTAA GTTGATAAGC TGTTTTTAAT ACAGTAGCA TCTCTTCAA AGTCGTCCGC	4200
TGAACACCA CAAGACGCTT AAATCGTGA TCAGTTAGTT GTTACTTGC TTCATAATTC	4260
ATAGAACTAT AGTAAATGA AATAAGAA CAAGATAAATCG ATCAGGACAG TCAAAATCGAT	4320
TCTCAACAAT GTTTTGAAG TAGAGGCGTA CTATCTTAGT TTTCAATCTAC TATACTATAC	4380
CATATTTTGT TTGCAAGGA ATCTATTATA AAAGGTTAAG TATTGCAAAA ACACTTAACC	4440
TTTCTTTTAA TACTTCATTA AGCTCTACTT TTATATAATAC TTCAAGCCCC ACATGAGCAG	4500
AAGCATGATG ATTAAGCAGA GAACAGCGCC AATATAAGCG ATTATTGTT GGTAGGATTC	4560
TCTGCTGTG ATACCTCTAT ACAACAAAAT AATAGACATA AAACCTGTCA AGCAGATGAA	4620
CATAAGTTGA TTGCTTCAG GACTAACCAA ATCATCATCT TCAAACTCTC TTATCTCTCAT	4680
TTCCCTATG AGATAAACAG TAAACAAAAT AGAAGCCAG TTAATACTA CTAAAGAGAA	4740
TTGGAAACT ACGGAAAAAT TTAAAACTG ACGAGATAGA AATAGATAAG TAGAAACAG	4800
CAAGGGCAAC TGACCTAAGA ACAATCTCGC AAGGAAGATG TTCCGTTTTT TAGCAAGAAA	4860
AGTTTTCATT TCTTTCTCC TTCTTTTAA TTGATAGCA AATAGATCAT AACTGCAATC	4920
ACATAGGCTA TGGTATAAAA TAGCTGATAC CAAGCACTCT CCCTAAGCG ATATAGAAAG	4980
ATGGACATGA TTAGATACAG AACGAAAAA ATCAGTATTT TTTCTTCAT AAGATTTCTT	5040
CCTAAATGTG CGATTTATCT TAGTTGAGCA AGAACATTTA CACTGCTAGT ATAGCACTTA	5100
TTTTGACCTT GGATCACTCA AATCATAAAT GGTATCAAAA ACCTCTTGAA TTGTAAAAAT	5160
TAAAAAGCA AGCATGAAA ACATACCTTC CTCTTTTAT TGTATTGATA CCAACTTGTT	5220
TGTAGACTTT TCATCTGCT ATCACAATC ATTTTGACAG GCGAAACAA ATTAAGAGAA	5280
CTCCCCGTA AATTAAAGCTA GCAAAATCAG GGGAGAAAT TATTTTATTAG AGAGTACTAT	5340
CCGTATCCTT TTGGAAGAT TTGAAAAATA TTTTCTAAT TAAGTCATCC ATATAGGAC	5400
CAATATATACC AACTACTAAA CCAATAATAA AACTTTTAAA ATCCATAATT ACCACCAACA	5460
TGTTGCTGCA TAGGCTACAC CTCCAAGTAT AGCTCCACCC GCAGCACCAG TTGCTGCACC	5520
TTGCCATGTT CCTGTTTTAA TGCCTAGTTG AAGACCTCTT GCTGCTCCTC CTCCAACACC	5580

TGCTTTGGCA AAATCTCCCC AATTGCATCC GCCACCTTCA ACGCAAGCAA GCATTTCAGT	5640
ATCCATAACA GAAAATTGTG ACATCATTTT TGTATCCATG ACAAACTACTC CTTTTTTAAA	5700
AAACTAAAT AAATCAGAA AGAATCCTCA TAATTTCAT ATAAGTCTTA CCAACTTAGT	5760
CCCAATTAT CACCAACCAT ACCTCCTAAG CATGTTAATC CACCCCAAT TGCAACAAAG	5820
TGTGCTCCAA CAAATGCCAG AGCAAGTCCA GCTACTCCTA AAGTGGCCAA ACCTGCTCCA	5880
GTTCCACCAG TTATAATTCC CGTAGTGACT CCTGTAAATCA GTGCATTTTG ACAATCAGTG	5940
GAGCTATACC CCCCTTCAAC TTTCGCAAGC ATTTTCAGTAT CCATAACCTC TAACTGTGAC	6000
AACATTTTGT TATTTCATGAT GAATACCTCC TTTTATTTT CAATTGTGTA CCAAAGTCTT	6060
AAATTCAATA AACAAATAGA TTTTTATAG TATCTTTTG ATTTTCTTAA AAAAGTATAT	6120
ACGTCTACTA TCTTCTTAA GGTAGCAGTA CCTATTTTT ACTCTAAGAT TTCAATAATC	6180
TTGAGTATCT AAAATATCTT AATTTCGTTA TTCTCCTTGC AATAAAAAGT TTTACTATAC	6240
TATTTATTA CTTCGAGAAA GCAAAAAATA TTAGTAAATA ATAGTTTATA GTTAAGTTTT	6300
TTATTCCTAC CAATCCATCA ACTAAGTAAA GCATCAACGA TTACATAAAC GATTGATAAT	6360
ATAATFAAAA TTTTGTCAAC TATCTTATTC TCATCATTTCT TAGATAACTT TGATATTTTG	6420
TAAGTAAGTA AATAAGACAG TAAATTAATA GCGATAATA TACTATATTT AAGAATCATTA	6480
ATCTTACAAA GAGGACATAA TTCTGTAACC TACACAATA AGTGTGCTG CTCCCCAGT	6540
TATCGGACCA GTCGCAGCAG CTAATAGTAC TGCTCCAATA CAACCACGA TTGCAGATCC	6600
TAAATTGCCT CTTCCTCCAC TAATATTTTC GAGTCTTCCA TTATCCATAA CAGAAAAATTG	6660
TTCCATCAT TTGTGATCCA TGACAAATAC TCCCTTTTTC TTTTATTTAT TTGTCTTGT	6720
TGTAACCTTG ATAAGTTAG TATATCATCG TTTTAAAAA TTTTTCATCC AGATCTTGAA	6780
TTGTCATCGA AACGCTTGA ATTAGCTTTT TTATTTCAAG CCACCTCTAA ATGTTTAAAA	6840
AAAAATATTT CTAATCACTT TTTTACCAT CAGGAAGTTT TAATGACTAT TCAAGATTTC	6900
ATAAATATG AACTTAGTTT TATGACATAA TAGACCTATC CACTATATGA AAGGAATTGC	6960
CAATGACTTC TTATAAAGT ACATTTGTTT CTCAAATAGA TGCGAGAGAC TGTGTGTCTG	7020
CTGCTCTAGC CTCGATGCTT AAATCTTATG GTTCAGATTT TTCTCTAGCT CACTTGAGAG	7080
AACTTGCAAA GACCAATAAA GAAGGGACGA CTGCTCTTGG CATTTGAAAA GCCGCTGATG	7140
AAATGGGCTT TGAACAAGA CCTGTTCAAG CAGATAAAAC GCTCTTGAAC ATGAGTGATG	7200
TCCCCATACC ATTTATCGTT CACGTTAACA AAGAAGGAAA ACTCCAACAT TACTATGTTG	7260
TCTATCAAC AAAGAAAGAC TATCTGATTA TTGTTGATCC TGACCTTCT GTAAATATCA	7320
CTAAATCTC AAAAGAAGC TTTTCTTATG AATGACTGAG AGTAGCTATT TTTCTAGCTA	7380

CCTCTGATTT	TCAAACAAAA	ATCTCTCATT	GCTTACAATG	TCTCTCAAG	CTTATGGTC	7500
ACTATATATCA	ATATAGGTGG	TCTTACTAT	CTCCAAGGAA	TCTTGGATGA	ATACATTCCA	7560
AATCAGATGA	AATCAACTTT	AGGAATCATC	TCAGTTGGTC	TGGTATCAC	CTATATCCTC	7620
CAACAGTCA	TGAGCTTCTC	CAGAGATTAT	CTCCTAACCG	TTCGAGTCA	GAGATTAACT	7680
ATTGATGTGA	TTTTATCCTA	TATTCGCCAT	ATTTTGAAC	TCCCATGTC	TTTCTTTGGG	7740
ACAGCTGCTA	CAGGAGAAAT	CATTTACAGA	TTACAGATG	CTAACTCTAT	TATAGATGCC	7800
TGGCTTCTA	CAATTCCTTC	TCTTTTCTG	GATGTTCTA	TTCGATTCT	TGTAGGAGCC	7860
GTCTTACTGG	CACAAAACCC	TAATCTCTTC	CTTCTTCTC	TATTTCCAT	TCCTATATAC	7920
ATGTTATCA	TCITTTCTTT	TATGAACCT	TTCCAAAAA	TGAACATGA	TGTCATGCAA	7980
AGTAATTCTA	TGGTATGCTC	TGCCATTATC	GAAGATATCA	ACGGATTGA	AACTATAAAG	8040
TGGCTACAGA	GTGAAGAAAA	TGGCTATCAA	AATATAGACA	GCGAATTTGT	AGATTATTGT	8100
GAAAAATCCT	TTAAGCTCAG	TAAATATTCT	ATTTTACAAA	CGAGTTTAAA	GCAGGAAACA	8160
AAATTAGTTC	TGAATATCCT	TATCCTATGG	TTTGGGCTC	AATTAGTCAT	GTCAAGTAAA	8220
ATTTCTATCG	GTGAGTGTAT	TACCTTTAAC	ACACTTTTTT	CTTACTTAC	AACTCCTATG	8280
GAAATATATTA	TCAACCTCCA	AACCAAACTC	CAATCTGCGA	AGGTGCTAA	TAACCGTTTG	8340
AACGAATCT	ATCTAGTCGA	ATCTGAATTT	CAAGTTCAAG	AAAACCCGT	TCATTCACAT	8400
TTTTTGATGG	GCGATATTGA	ATTTGATGAC	CTTTCTTATA	AGTATGGTTT	TGGATGAGAT	8460
ACCTTAACAG	ATATTAACT	CACGATTAAA	CAAGGAGATA	AGGTATGCT	AGTTGGAGTT	8520
AGTGGTCTG	GTAACAAAC	TTTAGCCAAA	ATGATTTGTA	ATTTCTTTGA	ACCTTACAAA	8580
GGGCATATTT	CCATCAATCA	TCAGGATATT	AAAAACATG	ATAAAAAAGT	CTTGGCCGT	8640
CATATTAAAT	ACCTACCCCA	ACAAGCCTAT	ATCTTTAATG	GCTCTATTTT	GGAAAACTTA	8700
ACCTTTGGCG	GTAATCATAT	GATTAGTCAA	GAAGATATTC	TAAAGCTTG	TGAAGTAGCT	8760
GAATCCGTC	AAGACATTGA	AAGAATGCCT	ATGGGCTATC	AAACTCAGCT	CTCTGATGGA	8820
GCTGGTCTAT	CAGGAGACA	GAAGCAACGA	ATCGCTCTCG	CTCGTCTCT	TTTAACTAAA	8880
TCTCTGTGTT	TAATACTAGA	TGAAGCTACT	AGCGGTCTTG	ATGCTCTGAC	TGAGAAAAAG	8940
GTTATAGATA	ATCTTATGTC	TCTAACTGAT	AAAAACATTC	TCCTTTGAGC	CCATGCTCTC	9000
AGTATAGCCG	AACGAACCAA	CCGTGTCTATT	GTCTTTGACC	AGGGAAAAAT	CATTGAAGTT	9060
GGTA						9064

250

(2) INFORMATION FOR SEQ ID NO: 18:

- (i) SEQUENCE CHARACTERISTICS:  
 (A) LENGTH: 7780 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 18:

CTCCATTTT TTGATTTTCAT AAATAAACAA CCTCTCTGTT AATTTGTAT ANTATAACG	60
ATATCCAAGT TACTTGTCAA GFGTTTTTTA AATTTTATC TCAAAATAT TTTTTCGTC	120
AAAAAAGGA GCCATCAGTT GATTTCAGC TCCCTTTTAT ACAGAAATTA ACTATTTTAT	180
AGTTGACAA TCTTACCTGT TTCAAAGTAG ACAACCCATT CACAGATATT TTTAGCATAG	240
TCACCGATAC GCTCCAAGTA GGAATAAAT TGAATAAAT CACGACCGT AACAAATGGCT	300
TCTGGATTTT TCTTAATCTC TTCAGTCGCA AGGTCACGGA TAGTTTCAA ATAGTGGTTA	360
ATTTGCTCAT CCATGGAGGC CACCGGTAT GCGTCCTCA CAGAACCAAT AAGATAAAGA	420
TCAAGTGTG CTTCACAAC GCTTTTAACT TCACGTCCA TTTTTTAAAT TTCTTCCTCT	480
ACAGCTGGAA TGGGCTCTC CCCCTTCATA CGGATGGTTG CTTGGGCAAT GGCTACAGCG	540
TGATCCCGCA TAGCTCCAC ATCTGATACA GCTTAAGGA CAGTCAAGAC TGTACGCAA	600
TCTTGAGAGA CTGGTTGTTG GAGTGGGATC ATTTCAATG ATTTCTTTTC CAGTTTCACT	660
TCGTATTCAT TTACTTCTCG ATCATCTTCG ATGACCTCTT TTGCCAGGTC ACGGTCATGC	720
GTGACAAAAG CAGTACCGT ACGATTGAAT TGTGAGAGCA CTCTTGTCC CATAGCOTAG	780
AATGTGTTAT GTAATTTCTC TAAATCTTCT TCAAAATTAG ATCGTAACAT CTTTCACTC	840
CTTATCCAAA TTTTCTCTGA ATATAGTCTT CCGTTTCCTT GTGTGGGGA TCAAGGAACA	900
TCTGCTTGGT ATCATTAAAT TCAATCAAAT TCCTATCTAG GAAAAATCCT GTCTTATCAG	960
AGATAGCTGA AGCTTGCTGC ATGGAACGGG TTACCAGAAG CATGGGTAC TTGCTTTTAA	1020
GACCATACAA GGTTCCTCA ATTTTACCAG CTGAAATCGG ATCCAAGCC GAAGTTGGCT	1080
CATCCAAGAG GATGATTTTA GGACTAGTTG CCAAGACACG GGCCACGAGC ACACGCTGCT	1140
GTTGACCACC TGACAAATCCA ATAGCTGAAT CATATAGACG ATCTTGACC TCATCCCGAC	1200
TAGAGCGACC TTGCAAGGCT TTTTCTACGG CTTCATCCAG AACCTGCTTA TCCTTAATTC	1260
CATTGATACG AAGCCCGTAG ACAACATCTC CATAGATAGT CATAGGGAAA GGATTAGGTT	1320
GTTTGAAAAC CATTCGGATT TCTTACGTA ATTCAAACCG ATCTGTACGC GGAAGTGTAG	1380
TGTTGTGACC ATTGTACACC ACGGATCCAG TGTGCTCAC CTCTGGATTG AGATCTCCCA	1440



TCGCGTTGAG	AGACTTGAGG	AGGGTTGACT	TCCCCTGATCC	AGATGGACCA	ATCAAGGCTG	1500
TAATTTCCCT	AGGTTGGAAA	GATAGGGAAA	CACATATTCAA	AGCCTTCCTT	TTATTATAAAT	1560
AAACGGACAG	GTCCTGATACC	TGTAAATCG	CATCTGTCAT	ACGGTTTCCT	TTCTAACCAA	1620
ASTGACCAGA	TACATAGTCA	TTGGTGGACT	GTAGCTTGGC	ATTTTGGAAA	ATAGTTGCAG	1680
TCCTGTCTATA	CTCAATCAAA	TCACCCAAGT	AAAAGAAAGCC	TGTATAGTCA	CTTGACGAG	1740
CAGCCTGCTG	CATATTATCG	GTTACAAATGA	TGATGGTAAA	GTTTTCCTTG	AGCTCAAAAC	1800
TGCTCTCTTC	TAGTTGCATG	GTCGCAATCG	GATCCAAGGC	TGAGGCTGGC	TCATCCATTA	1860
AGAGGATATC	TGGCTTAACA	GAGATGGCAC	GAGCGATACA	GAGACGTTGT	TGCTGACCAC	1920
CTGATAAGTT	CAGGCTGAC	TTGTGGAGAT	CGTCTTAAC	CTGATCCCAG	AGGGCAGCCT	1980
GACGAAGGGA	GTTTCTACG	ATTTCAATCTA	GGACTTGCTT	ATCCTTAACT	CCAGCACGTT	2040
CATGCGCAAA	GGAATATTA	CGGTAATTC	ACTTAGCAAA	TGGATTGGGA	GTTTGAAGAA	2100
CCATTCCAAT	GTTTACGC	ATTTCAATAA	CGTTGATTTT	TGGACGGTTG	ACATCAATTTC	2160
CACGATAGAG	AATCTGCCCA	GTTACTTTAG	CAATATCAAT	AGTATCATTC	ATGCGATTGA	2220
GACTGCGTAA	GTAGGTAGAT	TTCCCGGATC	CCGACGGGCC	AATCAAGCT	GTAATTTTAT	2280
TTCTTTCAAA	TTGCATATCA	ATCCCCTTAA	TGGATTCAAT	TTTACCATAG	TAAACATGGA	2340
CATCCTTAGT	AGAAAAGGCT	ACTTTTCTTT	CAGGAAAGGT	AAGGATATGC	TTCTCATCCC	2400
AGTTATATGT	TGACATGGCT	TCTCCTTTAG	GCAGCGGTTA	ATTTCTTGTG	TAGATAGCTT	2460
CCGAACCTTAC	GAGCTCCAAA	GTTAAAAATC	AGGATAAAGA	TCAGGAGCAC	AGCGGCAGAA	2520
CCTGCTGATA	CAATGTTTCC	ATCTGGAATA	GTGCTTCAC	TATTTGACTT	CCAGATATGG	2580
ACAGCAAGGG	TTTCTGCTTG	ACGGAAGATA	GAGATGGGGC	TAGTCACACT	GAGGATATTC	2640
CAGTTAGACC	AGTCAAGAGC	TGGGCGCGAT	TGCCCTGCTG	TATAGATCAG	AGCTGCAGCT	2700
TCGCCAAAGA	TACGACCAGA	TGCCAAGACG	ACACCGTTTA	CAATACCTGG	AAGCGCTTCC	2760
GGAATAACAA	CATGAACAC	TGCTTCCAG	CGAGAAATCT	CAAGAGCCAG	ACCAGCCTCA	2820
CGTTGGGTAT	GGTGAACGTG	TTTCAAACTA	TCTCTACAT	TACGCGTCAT	CTGAGGCAAG	2880
TTAAAGACTG	TCAAGGCCAA	GGCACCTGAA	ATGATTGAAA	ATCCATACTC	AAACTGGACT	2940
ACAAAGATCA	AGTAACCAA	GAGACCCACC	ACCACCTGATG	GTAAAGAGGA	CAAAATTTCA	3000
ATACAAGTCC	GCACAAAGTT	GTTAACAGGA	CCTTTTTTAG	CATATTACAG	CAAGTAAATC	3060
CCAGCTCCCA	TAGAAAGAGG	TACAGAAATA	ATCAAGGTAA	TGACCAATAG	GAAAAGGAA	3120
TTGTAAAGCT	GAATGCCAAT	CCACCAACCT	GCTTGAAAAG	CAGAAAGCCT	TCCAGTCCAG	3180

252						
AAAGACCAAG	AGATATGGGG	CAAGCCCCGA	ACCAAGATAT	AGAAGATCAA	GGAAGCCCAAG	3240
ATTGTCAACAA	TGATGCTAGC	AATCGTATAG	AGGACAGCTG	TTCGCAAGTTT	ATCTAATTTTC	3300
TTAGCGCGCA	TAATTTTTTCT	TTCTCTTTTC	TTTCGTAAATC	AATTTAATCA	CACTGTTTAAA	3360
AACTAAGCTC	ATCAAGAGCA	GTACCAAGGC	CAGTGACCAG	AGAACAATTA	TATTTACAGT	3420
TCCCATGACA	GTGTTCCCAA	TTCCCATAGT	TAAATAGAAA	GTAAAGTTG	CAGCTGGTGT	3480
GGTCAAGGAA	GTTGGATAA	CAGCTGAGTT	TCCGACAACC	ATCTGATAG	CTAGAGCCTC	3540
ACCAAGGCA	CGCGCATCC	CAAAGACAC	TGCAGTGAAC	ATACCAAGAC	GGCCCGCCTT	3600
CAAGATCACA	GCCCAAGATG	TCTGCCAGCG	AGTGGCTCCC	ATAGCGAACC	TGGCTTCACG	3660
ATAATAACGA	GGAAACCGAC	GCAAGCTATC	CGTTGTCTATA	AAGGTTACGG	TGCGCAAAAT	3720
CATGACAAAG	AGGACGGAAA	TCCCTGACAA	AATCCCAAAA	CCAGTCCCA	CAAAGACACT	3780
GCGAACAAAG	GGAAACGAGA	CTTGCAGGCC	AATAAATCCG	TACACTACTG	AAGGAATCCC	3840
AACCAGAGCT	TCAATAGCTG	GTTGCAAAAT	CTTCGCCCTT	TTTGGTGATA	CTTCGGTCAT	3900
AAAAACTGCT	GCACCAATAG	CAAAGGGTGT	TGCGATAAAG	GCTGAGAGAA	TGTTAAGCAT	3960
AAAGGAACCC	AAAATCATAG	GAAAGGCACC	AAATCTTTTA	CTAGAAGGAT	TCCAAGTTCC	4020
TCCCAAAAGA	AAGTCAAAGA	TATTCACACC	ATTGACAAAG	AGGTCGACA	AGCCTTTTTG	4080
CGCTACGAAA	ACCAAAATCA	TGGCCACAAG	GATGACTATC	AAAGAAAGAC	AGGCAAGGTT	4140
CAAACCTTTT	CCTAATTTCT	CCAGACGAGA	ATTTCTTGAT	GGAAGCAACA	TTTTCTTAGC	4200
TAATTTCTCT	TGATTCATTA	TTGTCTCCCT	TCCAACACTG	TCACAGTTCC	GGCAOCATCT	4260
TTTTCAACCT	TCAATTCCCT	AATCGGAATA	TACTTCAATC	CTTTGACAAAT	CCCTCTCTGG	4320
GTCTCATCCG	AGAGAACAAA	ATTGAGAAAT	TCTGCACCA	ACTCATTTGG	CTGCCCAAT	4380
GTATACATAT	GCTCATTAAGA	CCACAAGGCG	CAATTATTGC	TACTTATATT	TTCGTGACTT	4440
AAGTCATAGC	CATTCAACTT	CATGCTTTTG	ACCGAATCAT	CTATATAGGT	AAGAGATAAA	4500
TAAGAGATAG	CTCCTGGACT	TTTTGATACG	ATTGATTTTA	CCGCTCCATT	TGAATCCTGC	4560
TCTTGACTTT	GCATCGGACA	CTGACCTTCC	ATAATGACAG	TATCAAAGGT	AGCAGGAGAG	4620
CCAGAGCCCG	CTGCCCGATT	GATAACAGAG	ATGGGTAAAT	CCTTACCACC	AACTCTTTTC	4680
CAATTTGGTTA	CCTCACCTAT	GAAGATTTGA	CGAAGTTGCT	CTGTGGTTAG	GTATACACAA	4740
TCAACCTCCT	TATTTGACAA	CAGAGCCCAAG	CCAGCTACCG	CGACCTTTTG	GTCAACAAGA	4800
GCAGAAGCAT	CAATTCCTGC	TTTTTCTCTCA	GCAAAATACAT	CTGAGTTTCC	TATATCAACT	4860
GCCCCGAGCT	GAACCTGGGA	CAAGCCTGTA	CCAGAACCTC	CCCCTTGGAC	ATTGACCGTT	4920
TTTCCCAACAT	GGATCTGTGC	AAATTCATCT	GGCGCTACTT	CAACCAAGGG	TTGCAAGGCA	4980

GTGTAGCCAA CAGCGGTAT GGATTCTCCA CGATCAATCC AGCTAGCACA GCCTACTAAA	5040
CAAGCCGTCA GCCAAAAGC GATAAGAGAC AGAGCAAGCT TTTTCTCTTT TTTCACGTGT	5100
TTTCTCTCG AAAATAATTA TGAATACTGT GAATTTTTTA AGTAGTTCTT TATGAGTTGA	5160
CGCATGAATT CTTACCAAAAT TTCTGCGCAA TTGATTATTT ATATAATATA GGCTATATTA	5220
CTCTTTCTTA ACCTCCTTTT TTCAATATGT GATAAAATCT CTGTCTATC CCTTCCCCCA	5280
TTGTACCCCA TTATAGTCAT TTCTGTCTCT TTTTCCCCCT TTTTAATGCA AGGGAATTA	5340
CTCTCCTTAG ATGATAATCC AAAAGCTAGA AAGGTATCTC AAACCTCTCT ACTCTCCCAG	5400
ACTAGTTTAC AACTAAAAGG AAAAGATTCT ATTTTATGAG AAATCTAGTT TACAAGCGGT	5460
AAGAACGCTA ATAACTAAAC TTCTTGTACT CTTTGAAAT CTCTTCAAC CAGTGTTTTG	5520
AGCTATCTAT GGCTAGCTTC CTAGTTTGCT CTTTGATTTT CATTTGAGTAG TAAAACTACA	5580
TGTAATGCA ATCAAGATAT CAAGAATCAT CCTACTAAAA AAATCCATAC TTTCACTATA	5640
ACATAGATA AGATATTGTA CTAGCATTTT CATTTGAATC TGAGGCCTTT TGGAAAAATA	5700
TTTTTCAAAA CATTTCCAGT AACCTTTGCA AAGCCCAAG CATTTGCCCTT AACCAAACT	5760
TGGTACCAAC CATTTGGCAG ACTTCTTGCC AGCTGAACGG TTCTCCAGC CGCATACTTG	5820
ACAAAGCTT CTTGGCCAAT TTCAACCGAC TGTTCGACCT GACTCGGTTT CAAGCTAAAA	5880
CCAGAGCGA AACTGGGCTC AAAGCGTTTC TTCTTAAAG TACCCAGATG CAGTCCATTG	5940
CGAGCAATCT TGAGCTTCCA TAAATCTGGC AAAAGTCTG GCAAGAGATA AAGCTGGTCT	6000
CCAAAAATCT GCAAGATACC CGGTAGATTG ACCTTCAAAAT GGTTTGGGC AAATTCCTGC	6060
CACAAGGCAA CTTTGTTACG GCTGAGGTTA CTCTTACTTG CCTTAAATTT AGGAGCTGGA	6120
TTGTTACCTT TAAACTGTAG ATGGGCAACA AACTGACCTT CTCCCTTAAA CTGATGAGGA	6180
TACATCCGAG CCGTTTCTGG CAGGTCAATA CCAGTACCA TTCCATTGAT ATGCTCTACT	6240
GGCAACAAAT CAAAATCATA CTCTTCCAGC AACCAATTGA CAATCTCTTC GTTTTCTCTG	6300
GGTGCCAGG TACAGGTGCA ATAAACCAGA TGACCACTTT CAGCTAACAT GGTCACTGCA	6360
TCTTCCAGAA TTTCTCTTTG CAAGCTAGCA CATTTGACTG GATAATCTAA GCTCCAATAG	6420
TCCATAGCAT CAGGTGCTT ACGAAACATT CCTTACCAG AGCAAGGGC ATCAAGAAGC	6480
ATTAAAGTCAA AATGACCTTT AAAGACCTTG ACCAAGCGGT CGGCAGATTC ATTGGTCACC	6540
ACGACATTTG TCGCTCCAAA ACGTCCCATG TTTTCAACCA AAATCTTACG CCGTTTGCTT	6600
GAATTTTCAT TTGGAACAAG TAGCCCCCTC CTGTCTAGAT AGGCTCCAG TTGAGTTGAT	6660
TTGCCCCCCG GTGCAGCAGC CAAGTCCAAG ACCTTCATAC CAGGACTGGG TTGGCTACT	6720

254

TGAGCCACCA TTTGAGCAGC AGGTTCCTGC GAATAAACA AACCCTGAGC ATGCTCAGGC	6780
GATTTCCTCG AAACCTTCCTC ATAGTGGCCC CAAGGGGTTT GAGTAATGGC ATCAGAAAG	6840
GAAAGTTGCT CTCTCTTTAA GGGATTGACC CGAAAGGCCG AAACCGCTTC CTCTCAAAA	6900
GAGGCAAGAA AATCTCTGCG CTCATCTCCT AGTATCTCTT TATATTTTTT AACAAATCCT	6960
TCTGGAAATT GCATTTAAGT TCTTTTCCTT TGTAAATAT AGGACTGAAT TTCTCTCTGC	7020
ATCTCAAGAG GCACCATCAT GACCGGCTGT CTGTTTGAA AATCAGGAGC TTCACCAAAA	7080
AGGGTCACAA CCGATAGCC CAGACTTCC CTTAAAATAC TAGCTCGCGC ATAATCCCAT	7140
GTTTGCAGAT AAGTGAGATA GGTCAACAAA CGCCCTGACA AAATCTTGGC AAAACTAATG	7200
GCCTCACTTC CATAGACAGC AACACCAAGA ACCGCTCGGC TCAATCAGC CAGCCCCCAT	7260
TCATTGGTTT CCAGCATACC ACTATTCCCT GCAATGAGAA AATCTCCAAG TGGTTTAGTT	7320
TTAAAGGAG CTAGGACCT ATCATTAGA CAAACTGAA ATTCGCCACC ACCGTGGTAA	7380
CAATCCCCCT TGACCACATC ATAAATCAGA CCAACTGTC CCTGACCATT TTCAAAATA	7440
GCCATCATAA CAGCAAAATC TTCTGCTGG GCTACAAAT TATTGGTACC ATCAATGGGA	7500
TCAATGACCC AAACCTTGCC CTCTTGAACC GAGGCTCGCA GACAACTTC TTCAGCACAA	7560
ATCTTATCCT CAGGATACG GGACAAAATC TCACCAACCA AGAGTTCTCG AACTTCTTTG	7620
TCCAGTCTGG TCACCAATC TGTGGAGAG GACTTGGTTT CAACACGAA GTCTTCTGTC	7680
ATATGGTCAA GAATGACTG ACCTGCTTTC TTAACAAGCT CTTTAGCAA TTCAAATTTA	7740
CTTTCCAAGA GAAATCTTTC CTTCCTCTTT TTCTTTGGGG	7780

(2) INFORMATION FOR SEQ ID NO: 19:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 4820 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(x1) SEQUENCE DESCRIPTION: SEQ ID NO: 19:

GTAATGATAT AGGAACACCA GGTGACCTGA TGGGACGTCG TAAGCCTATG AACTACTAGC	60
TGCTAAAGGC TTAAAGATG GTATGGTACC ATATATCTCA AACCAATACG AAGAGAGAGC	120
CAAAACAAAG GGCAGACAA TCAATCTCTA CGGTAAACA AGAGGTTTGG TTACAGATGA	180
CTTGGTTTGG GAAAAGGTAT TTAATAACCA ATATCATACT TGGAGTGAGT TTAAGAAAGC	240
TATGTATCAA GAACGACAAG ATCAGTTTGA TAGATTGAAC AAAGTTACTT TTAATGATAC	300
AACACAGCCT TGCACAAATC TTGCCAAGAA AACTACAAGC AGGTAGATG AATTACAGAA	360

ATTAATGGAC GTTGCTGTC GTAAGGATGC AGAACCAAT TACTACCATT GGAATAACTA	420
CAATCCAGAC ATAGATAGTG AAGTCCACAA GCTCAAGAGA GCAATCTTTA AAGCCTATCT	480
TGACCAAAACA AATGATTTTA GAAGTTCAAT TTTTGAGAAAT AAAAAATAGT GTCTACTATT	540
AGGAAATPAA GTTTAAAAAG GTGATGAAGA ACAAAACCAAG ATTCAAGCAG GAATTCCTAC	600
TGATAATGAA GTAAGTTATG ATCTTATTTA TCAGCAGGAA ACTCTTCTCG CAACAGGTTT	660
ATCAACTTCT GAGCTTACAG CTTTAGGCCCT ATTAGCTGTT GGTAAGTTAG TTCTTTTGGT	720
TCATAATATG ACGGGAACAG TTTTTTGCTC CCTCTGAAAA GTCATCATTT GATGGCTTTT	780
TTCTATATAG GGTAAAAGAT AGGGTAAAG GCTATCATCG GACAAAATAA AGAAGGCATG	840
ATATAATATA AAGTAGATTT CTATGTGATA AAACAAGAAC TGTTTGGACA TCATTCAATT	900
GAAAACTCTC TATGTTCAAA CAATAGTAAA ATAAAAATAGG GGATCTAAAT CCTTGCTATG	960
AJAGGAAAAA ACTCAATGGC TACTATTCAA TGGTTTCTCG GTCACATGTC TAAAGCTCGT	1020
CGACAGGTGC AGGAGAATTT AAAATTGTTT GATTTTGTGA CGATTTTAGT AGATGCACGC	1080
TTGCCCTCAT CTAGTCAAAA TCCTATGTTG ACCAAGATTG TTGGTGATAA ACCAAAACTC	1140
TTGATTTTAA ACAAGGCCGA CTGGCTGAT CCAGCAATGA CCAAGGAATG GCGTCAGTAT	1200
TTTGAATCAC AAGGAATCCA GACGCTAGCT ATCAACTCCA AAGAGCAAGT GACTGTAAAA	1260
GTGTGAACAG ATGCGGCCAA GAAGCTCATG GCTGATAAGA TTGCTCGCCA GAAAGAACGT	1320
GGGATTCAGA TTGAAACCTT GCGTACTATG ATTATCGGGA TTCCAACGC TGTTAAATCA	1380
ACTCTGATGA ACCGTTTGGC TGGTAAAAAG ATTGCTGTTG TTGGAACAAA GCCAGGGGTC	1440
ACAAAAGGTC AACAAATGGT TAAACCAAT AAAGACCTGG AAATCTTGA TACACCGGG	1500
ATTCTCTGGC CTAAGTTTGA GGATGAACT GTTGCACTTA AGTTGGCATT GACTGGAGCT	1560
ATCAAGACC AGTTGCTTCC TATGGATGAG GTTACCATTT TTGGTATCAA TTATTTCAAA	1620
GAACATTATC CAGAAAAGCT GGCTGAACGC TPCAAACAAA TCAAAATTTGA AGAAGAAGCG	1680
CCTGTGATTA TTATGATATG GACCCGGGCC CTCGGTTTCC GTGATGACTA TGACCGTTTT	1740
TACAGTCTCT TCGTAAGGA AGTCCGTGAT GGCAACTCG GTAACATATC CTTAGATACA	1800
TTGGAAGACC TCGATGGCAA CGATTAAAGA AATCAAGAA TTCTTTGTGA CAGTCAAGGA	1860
GTTAGAAAGC CCTATTTTTT TAGAGCTTGA AAAGGATAAT CGCTCAGGAG TTCAAAAGGA	1920
AATCAGCAAG CGTAAAAGAG CCATTCAACG TGAATTAGAT GAAAAATTGC GCTTGAATC	1980
CATGCTTTCT TATGAAAAAG AACTTTATAA GCAAGGATTG ACCTTAATTG CAGGTATTGA	2040
TGAGGTGGT CGTGGTCTCT TTGCTGGTCC TGTAGTCCCT GCGGCCGTTA TTTTATCTAA	2100

256	
AAAATTGTAAG ATTAAGGTC TCAACGACAG CAAGAAATTT CCTAAAAAGA AACATCTGGA	2160
GATTTTCCAA GCGGTTCAAG ACCAAGCCTT GTCGATTGGA ATTTGGTATCA TAGATAATCA	2220
GGTCATCGAC CAAGTCAACA TCTATGAAGC AACCAACTA GCCATGCAAG AAGCAATCTC	2280
CCAGCTCAGC CCTCAACCAG AGCACCTTTT GATTGATGCC ATGAAACTGG ACTTGCCCAT	2340
TTCAAAAC TCATTATCA AAGGAGATGC CAACTCCCTC TCTATCGCAG CAGCATCTAT	2400
AGTAGCCAAG GTAACACGTG ATGAATTGCT GAAAGAATAC GATCAGCAGT TCCCTGGCTA	2460
TGATTTCGCT ACTAATGCGAG GATATGGCAC AGCTAAACAT CTGGAAGGCC TCACAAAAC	2520
AGGAGTTACC CCAATTCACC GAACCAGCTT TGAACCCGTT AAATCACTGG TTTTAGGTAA	2580
AAAGAAAGT TAAATTGAAG GAAATAACAT GGAGGAACAG TCGGAATAG TCGTTCCTAA	2640
GAAAGAAATC GCCTTTGCAT CCAGCACTAT ACTATCCCAA GTTGGTCGAG GAATCATTTG	2700
CGGCTCATC GTTGGAAATA TCGTCGGATC CTTCGTTTC TTAATTGAAA AGGGCTTCCA	2760
CTGATACAA GGAGTTTATC AAGATCAAGG GTACTTAGTG CGCAATCTTT TTGTACTGGT	2820
TTTGTTTTAT ATACTCATCT GTTGGCTCAG TGCCAACTA ACACGGTCAG AAAAAGATAT	2880
TAAAGGCTCA GGAATTCCTC AAGTCGAAGC CGAAGTGAA GGCCTCATGT CCTCAACTG	2940
GTGGGGCATT CTTTGGAAAA AATATGTGCT AGGTATTTCT GCTATTGCCA GTGGACTCAT	3000
GCTGGGTCCA GAGGGACCCA GCATTCAACT TGGAGCAGTT GGTGTATAAG GAATTGCCAA	3060
GTGGTCAAA TCCAGTCCAG TAGAGGAACG TTCCCTGATT GCCAGTGAG CTGCAGCAGG	3120
TTTAGCCGCA GCCTTTAATG CTCCTATTGC AGCACCTCTC TTTGTTGTAG AAGAAGCTTA	3180
TCACCATTTT TCGCGCTTTT TCTGGGTCTC AACTCTAGCA GCCAGCATCG TAGCAAACTT	3240
TGTGTCTCTA CTCATGTTGG GTTTGACACC AGTATTGGAT ATGCCAGATA ACAATCTCTCC	3300
CATGACCTTA GATCAGTATT GGATATATCT CGTCATGGGA ATTTTCCTTG GATTTTCAGG	3360
TTTCTCTAT GAGAAAGCTG TATTAAACGT TGGAAAGATT TATGACTTGA TTGGTCAAAA	3420
AATCCATTGG GATAGGGCTT ATTATCCCAT CTGGGCTTTT ATCCTTATCA TACCAGTCGG	3480
AATCTCTTTA CCTCAAATCA TTGGTGGCGG AAATCAGCTT GTCCCTTTCTT TAACGTGAACA	3540
AAATTTTAGT TTCCAAGTTT TATTAGCTTA CTTTTAAATC CGCTTTATTT GGAGTATGAT	3600
TAGCTATGGA AGTGGACTGC CAGGAGGAAT TTTCTCTCCC ATTTTAGCTC TTGGTCTCTT	3660
GCTTGTGCGC TTAGTTGGTG TTATCTGTGT CAATCTGGA CTGTGTAGTC AAGAGCAATT	3720
CCCTATATTT GTCATCTAGG GAATGAGTGG CTATTTTGGG GCCATATCAA AAGCTCCCTT	3780
AACCGCTATG ATCTCGTAA CTGAGATGGT AGGAGATATT CGCAACCTTA TGCCACTTGG	3840
CTCTGTCACT CTGTCTTCTT ATATTATCAT GGATTTGCTC AAAGGTACCC CAGCTATGTA	3900

257

AGCCATGCTG	GAAAAAATGC	TTCCAGAAGA	AGTATCTAGC	GAAGGAGAAG	TTACACTTAT	3960
CGAAATACCA	GTTCCTGATA	AAATPGCTGG	GAAACAAAGT	CATGAACCTCA	ACTTACCACA	4020
CAACGTCCTC	ATCACAACCTC	AAGTCCATAA	TGGCAAGAGC	CAACAGCTTA	ACCGCTCAAC	4080
CAGATCTGAT	CTGGGTGATA	TGATTCACCT	GGTTATTCCA	AAAAGTGAAA	TGGAAAAAGT	4140
CAAGAGTTTG	TTGTTGTAGT	ATGAGTATTT	ACATAAATTTA	TGTTATGTAA	ATGATCAGTT	4200
TGATTTTATTT	AGAAAAACGA	TTCTCAGGAA	TGAGTCCGGT	TATTTTCTAC	TGATGAGGAA	4260
TTTTACATAT	AAATAATTGA	ACTTTATTAA	AAATTAAGACT	ATAATTAAAT	TAGAAATGAT	4320
AAAGTATAAA	GCTAGAAAGG	AGTTTACTGT	ATCAAACTCG	TACAGTAAGA	TTAAAAATCAT	4380
GA AAAAGAAA	ACAATAGCAA	TTATATAGAG	AAATGAAATA	GAAATAGGAT	AAAAACAATCA	4440
GGACAATCAA	ATCAATTTCT	AGCAATGTTT	TAGAACTCCA	GATGTACTAT	TCTAGTTTCA	4500
ATCTATTATA	CAATGTGTTT	TGTATCTCAT	AGCTCCTTAT	ATAGCTCTTC	AGTTATGTAG	4560
TATTACAGAA	AGTTTAGTGG	GTGAGATTTT	TATTAATTTC	CTTATTTCTGT	TTTGTTTGTG	4620
GGCTTAAGTC	TTTTTATCAC	TTTGAAAAAC	TCTTATAACA	TCTTTCCGAA	AAACTATAAT	4680
TTTCTTGAAA	AATATACAAG	TCTATGCTAT	ACTACTAGTA	TACTTACTTA	TGGAGAAAAAT	4740
ACATGAAACG	TGAGATTTTA	CTGGAACGAA	TCGACAAACT	AAAAACAATC	ATGCCCTTGGT	4800
AACTTCTGGA	ATACTACCAA					4820

## (2) INFORMATION FOR SEQ ID NO: 20:

- (i) SEQUENCE CHARACTERISTICS:  
 (A) LENGTH: 21338 base pairs  
 (B) TYPE: nucleic acid  
 (C) STRANDEDNESS: double  
 (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 20:

CTACGACATC	ATGATTAAAC	GTATGCGCT	ACTACCAACT	GAGCTATGGC	GGATAAAATA	60
GTCCGTACGG	GATTCGAACC	CGTGTACCG	CCGTGAAAAG	GCGGTGCTTT	AACCCCTTGA	120
CCAACGGACC	TTCTATCTGT	AGCAGATATA	ACCATATATAT	CAATTCTTTG	CTAATTGTCA	180
ATCAGTTTGG	AGATTTTCTC	TCTAAAAATAT	CTTTTAATTT	TCTAATTTTT	AACTTTGAAA	240
TAGGACACAG	ATGGCTTTCA	TAGAAAAACAA	TTTCTAAGTT	TTTTCGATCA	ATTTCCTCTGA	300
TATTTACCTAT	ATTTACAAA	AATGACTTGT	GAGGAGAATA	AAATCGCTGA	GTATGTTTGT	360
CCTTTTCTCTG	AAATATCTGTC	ATGGTACCAT	AAAACTCTTT	TGCAAAATTC	TTACCAATAA	420

258

TGCCCAATTT ATGAGATACC CCTGTTGTTT CAATATACAA AATATCATGG TAAGGAATTT	480
TTAAATCATTT TCCCTTGTAAG TTGTAGTCGA AATAATCTAC AACATCTTCA TTTTCAAGTA	540
ACATACTCTT CGTGTAGAAG ATATTTTGCT CAATTCTCTT CTAAACATC TCATCATTGA	600
TATCCTTATC AACAAAATCT AGGGCTGATA CCTGGTATTT ATAGGTTAGA GTCCGAAACT	660
CTGATCGACT AGTGATAAAG ACGATAATAG CGTAAGGATT GTAATGACGA ATGAGCTGAG	720
CCACTTCAAA TCCCTTTTTC TCAATTCCAT GAAATTCGAT ATCTAGGAAA TAAAGCTGAT	780
TTACTTCAATC ATTTTCAATG TATTCTTCAA ATTCAAGGAC TTTTCCCGTT GTCTTGATG	840
ATATTGGAAT ATTCGATTCT TTCGAAATTT CATCCAATAT TCTCTCTAGT CTCACCTGAT	900
GTTCAATAAC ATCTTCTAAA ATTAAAACTT TCATTCAAAT TCCCTCTTAA ATCTAATGAT	960
TGTCTAAAT GTACTGCCTT CCATCTCTGT TTCTAAAAA ATATTGTTGT ACTTATCTAG	1020
TAGTTCCTTC ACATTATTTA ATCCGACGCC GCGATTTCCT CCCTTAGTGG AGAATCCTAA	1080
GGCAATAGA TCTCCTGAAG GAGTCATCGT CATTTTACAT GAATTCGAA TCACAATAAC	1140
TGTTTCAGTT TCCATCTTAA TAACTGCTAC TTCCATCTGC TTTTTATAGC TATCAGCCGA	1200
TCCTTCGACA GCATTATTCA ATAAAAAGCT CATGATACGA ACCAAATCCA ATAGTTCAAT	1260
TGGAAGCTTG GTAATCGTAT CTTTACTTTC CAGGTAAAC TCTACACCAT TATTTCGAGC	1320
ATAGACAATT GACTGAGCAA CCAAACTTCG TAAAGCTGAG TCTTCTATGT TGTTCAAATC	1380
AAAGTAAGTG TACTTATCTG AACGCAATTT ATGATTGTCT TTGACTAAAA CTTCAATTGT	1440
AATTCGTGCA ATTTCTCGTA AATTACCACT GTCAATTGCC ATCTGCATGC TGACAAGCAT	1500
TCCAGCATAA TCATGTCGAA AACCAAGGAT TTCAITATAC AGACCAACAA TTTCACTCTGT	1560
GTAATCTGT AAATGTTTCT GTTCAAAATTT CTCTGCTTC AAAGCAATCT CTTTCTCCAT	1620
TTGAATCTTA TGAGAATTCA TTGCAAAAGAA GGTCAAAAGG AGAGAGATAA AGACAATAGA	1680
TGACAAAATA CTTCCAAAAC TATTCAAATG TTAAATCGTA CTTACCATAT CTGAACGAA	1740
AGATACAATA TGTAGCAATA GTAAAGCAAA AATACTTTT TTCAAGAAAG GATAAAGGTA	1800
GTCTCTGTCA AAATAGGCTA GTTCCAAATG GAAATAGTAA ATGATTTTTA ATGTAACAAA	1860
ATAGGTTAAC ACCGTACAAA CGAAAAAGAA TGGGAAATGA TATTGTAAAA CAAAAATTATC	1920
TCTGTGTATA GAGGAGAAAA TTACGGACAG AAGGTTATGA GTGCTCTCAT ATAAAAAGAA	1980
TAGTAGTAAA CTTAGGAATA GTCTCTATC CCTCTCATC GTTTTCATCC ATCGAAAAATA	2040
GGAATATAAG CCCAAGGAA ATAAAAATCT TTCAATCCCT ATTTTATCTA AATATAGAAG	2100
ATAAAGGAA AATTCAAGTA CTATTTCAGT TAGTAATGTA TAAGCACCAA AAACGTATAA	2160
TTCTTTTCTA TTTATTCGAC CTTTACAAAT TAAACGGTAA CTGTGACTAA TAATTAAAAA	2220



ATGAACAATA	ACTGTCCCAA	ATCUAAGTAA	ATCCATTACT	CTTTCTCCTT	ATTTCATTAC	2280
TTTTTTTGGTA	GGAAAAGAAA	ATCAAGGATG	ATTCTTGAAA	TCCTCATCTC	CCCACCTTTA	2340
ATCTTTTGTGA	AGTCTTTTTC	CTTCAAAGCT	ACAAACTGTT	CCAAATTTAAC	TGTGTTTTTC	2400
ATATATAAAAT	CTCCTAAAAT	GTTTTTCTTT	GTAAGCTAAC	TTACAAAAAC	CATTATACAA	2460
AAATGGAATTT	CGTTTTAGAT	AAAATTTCTCT	CAACTGTCTAT	TTTCTCTCTC	CAAGGTGTAC	2520
TTTTTTAAGA	AAAAAGCCGG	GAAAATTTCCC	AGCTTTGCTA	TTATATTGAT	CCCAGCAGGA	2580
TTCGAACCTG	CGACCCCTTG	CTTAGAAGGC	GAATGCTCTA	TCCAGCTGAG	CTATGAGACC	2640
TAATACAAAT	ATTCTACCAA	AAATTCAAAT	AAAAGTCAAT	TTTCTATTTA	TGGTAGGGGA	2700
ATCCCTCTGTG	AATCGTAAAA	GCGCGATAGA	TTTGTTCAC	AAGAACTAGT	CTCATTAAC	2760
GATGGGGTAA	GGTTAGGCGA	CCAAAAGTGA	CAGAAAGATT	GGCTCTATTT	TTTACAGATG	2820
ATGATAATCC	TAAACTTCCC	CCAATAATAA	AAGTAAGAGT	AGAAAACTCT	TTTATAGAAG	2880
TTTCTTCTAA	CTGCTTACTA	AATTTCTCTG	AGAAGAAAAGT	TTTCCCTTCA	ATGGCTAACA	2940
CAATAACGAA	ATCACGGTCA	GCAATTTTGT	ATAAAAATCT	CTGACCTTCT	ATTCTTAAAA	3000
TCTTTTGATT	TTCTGATPCA	CTGGCCCTAT	CTGGTGTTTT	TTCACTCTGAT	AACTCAATCA	3060
TTTCAAACTT	AGCAAACTCA	GAAATTCGTT	TTGAATACTC	TGCGATACCA	TCTTTTAAAT	3120
ACTTTTCTTT	CAGTTTCCCA	ACTGTTACAA	CTTTAAATTT	CATGACTCTA	TTCTAACATA	3180
TTCTCTATTT	TTTCACATCT	TATTCACAAA	ATAAAAAATA	GATTTCAATT	AAGAAAAATCA	3240
CAATTTTCAA	AGAGTTATCC	ACAGTTTGTG	TAAAATTTTT	GTGTTTAAAGT	TATAATTAAG	3300
CTAGTCAGTT	TATACTTTCA	GTAATTCAAA	CATATGGAGG	CAATATGAA	ACATCTAAAA	3360
ACATTTTACA	AAAAATGGTT	TCAATTATTA	GTCGTTATCG	TCATTAGCTT	TTTTAGTGGA	3420
GCCTTGGGTA	GTTTTTCAAT	AACTCAACTA	ACTCAAAAAA	GTAGTGTAAA	CAACTCTAAC	3480
AACAATAGTA	CTATTACACA	AATGCTCTAT	ANGAACGAAA	ATTCAACAAC	ACAGGCTGTT	3540
AACAAAGTAA	AAGATGCTGT	TGTTTCTGTT	ATTACTTATT	CGGCAACAG	ACAAAATAGC	3600
GTATTTGGCA	ATGATGATAC	TGACACAGAT	TCTCAGCGAA	TCTCTAGTGA	AGGATCTGGA	3660
GTTATTATATA	AAAGAATGA	TAAAGAAGCT	TACATCGTCA	CCAACAATCA	CGTTATTAAAT	3720
GGCGCAGCA	AAGTAGATAT	TCGATTGTCA	GATGGGACTA	AAGTACCTGG	AGAAATGTGC	3780
GGAGCTGACA	CTTCTCTGTA	TATTGCTGTC	GTCAAAATCT	CTTCAGAAAA	AOTGACAACA	3840
GTAGCTGAGT	TTGGTGATTC	TAGTAAGTTA	ACTGTACGAG	AAACTGCTAT	TGCCATCGGT	3900
AGCCGGTTAG	GTTCTGAATA	TGCAAAATAC	GTAACCTCAAG	GTATCGTATC	CAGTCTCAAT	3960

260						
AGAAATGTAT	CCTTAAATCT	GGAGATGGA	CAAGCTATT	CTACAAAGC	CATCCAACT	4020
GATACTGCTA	TTAACCAGG	TAACCTCTGC	GGCCCACTGA	TCAAATTTCA	AGGGCAGGTT	4080
ATCGGAATTA	CCTCAAGTAA	AATTGCTACA	AATGGAGGAA	CATCTGTAGA	AGGTCCTGGT	4140
TTCCGAATTC	CTGCAATGA	TGCTATCAAT	ATTATTGAAC	AGTTAGAAAA	AAACGGAAAA	4200
GTGACGCGTC	CAGCTTTGGG	AATCCAGATG	GTTAATTTAT	CTAATGTGAG	TACAAGCGAC	4260
ATCAGAAGAC	TCAATATTC	AAGTAATGTT	ACATCTGTG	TAATTTCTCG	TTGGTACAAA	4320
AGTAATATGC	CTGCCAATGG	TCACCTTGAA	AAATACGATG	TAATTACAAA	AGTAGATGAC	4380
AAAGAGATTG	CTTCATCAAC	AGACTTACAA	AGTGTCTCTT	ACAACCACTC	TATCGAGAC	4440
ACCATTAAAG	TAACCTACTA	TCGTAAACGG	AAAGAAGAAA	CTACCTCTAT	CAAACTTAAC	4500
AAGATTTCAG	GTGATTTAGA	ATCTTAATTG	ACATCTATGT	AAAGAAGCT	TTACATAAGA	4560
GAAAGATGCT	GTTAGTGTAG	AATCATGGAA	AAATTTGAAA	TGATTTCTAT	CACAGATATA	4620
CAAAAAATC	CCTATCAACC	CCGAAAAGAA	TTGTAGTAG	AAAACTAGA	TGAAGTAGCA	4680
CAGTCTATCA	AAGAAAATGG	GGTCAATCAA	CCGATTATTG	TTCTCAATC	TCCGTATTAT	4740
GGTTATGAAA	TCCTTCAGG	AGAGAGACGC	TATCGGCTT	CACTTTTAGC	TGGTCTACGG	4800
TCTATCCCAG	CTGTTGTAA	ACAGATTTC	GACCAAGAGA	TGATGTCCA	GTCCATTATT	4860
GAAATTTTAC	AGAGAGAAAA	TTTAAACCCA	ATAGAAGAG	CACGCGCTA	TGAATCTCTC	4920
GTAGAGAAAG	GATTCACCCA	TGCTGAAATT	GCAGATAAGA	TGGCAAGTC	TCGTCATAT	4980
ATCAGCAACT	CCATTCGTTT	ACTTTCTCTG	CCAGAACAGA	TTCTTTTACA	AGTAGAAAAAT	5040
GGCAAACTAT	CACAAGCCCA	TGCGGCTTC	CTAGTTGGGT	TAAATAAGGA	ACAACAAGAC	5100
TATTTCTTTC	AACGGATTAT	AGAAGAAGAT	ATTCTGTAA	GGAAATTAGA	AGCTCTTCTG	5160
ACAGAGAAAA	AACAAAAGAA	ACAGCAAAAA	ACTAATCATT	TCATACAAA	TGAAGAAAAA	5220
CAGTTAAGAA	AATCTACTCG	ATTAGATGTA	GAATTTAAAC	TATCTAAAA	AGACAGTGA	5280
AAAAATCATTA	TTTCTTTTTC	AAATCAAGAA	GAATATAGTA	GAATTATCAA	CAGCTTGAAA	5340
TAAGGCTGTT	CTTTTATTTT	TTTATCTCAC	AAGGTTATCC	ACTATGTTTT	TCGATAAAAA	5400
GCTTAATAAA	TCAATATTTT	CTTCTTTTAT	CCCCAACCTG	TGGATAAAGT	TTGGTAACAT	5460
TGTGGATTAT	TTTTTCACGC	TTGTGAAAA	TTCTTGCTAT	CTATGATAAA	ATATCTCTAG	5520
TATTAACCTT	TTAAATAGTA	AAGGAGGAGA	AAGGATTGAA	AGMAAAACAA	TTTTGGAATC	5580
GTATATTAGA	ATTTGACAAA	GAAAGACTGA	CTCGATCCAT	GTATGATTTC	TATGCTATTC	5640
AAGCTGAAC	CATCAAGGTA	GAGGAAAAAT	TTGCCACTAT	ATTTCTACCT	CGCTCTGAAA	5700
TGGAAATGGT	CTGGGAAAA	CAACTAAAAG	ATATTATTGT	AGTAGCTGGT	TTTGAATTTT	5760

ATGACGCTGA	AATAACTCCC	CACTATATTT	TCACCAAAACC	TCAAGATACG	ACTAGCTCAC	5820
AAGTGAAGA	AGCTACAAAT	TAACTCTTT	ATAACTATAG	TCCAAAGTTA	GPATCTATTC	5880
CTTATTCAGA	TACGGGATTA	AAAGAAAAGT	ATACCTTTGA	TAACTTTATP	CAAGGGGATG	5940
GAAATGTTTG	GGCTGTATCA	GCCGCTTTAG	CTGTCTCTGA	AGATTGGCT	CTGACCTATA	6000
ACCCTCTTTT	TATCTATGGA	GGACCAAGCC	TGGGAAGAC	TCACTTATTA	AACGCTATTG	6060
GAAATGAAAT	TCTAAAAAAT	ATTCTCTAATG	CGCGTGTTAA	ATATATCCCT	GCCGAAGACT	6120
TTATTAAATGA	CTTTCTTGAT	CACCTAAGAC	TGGGGAAAT	GGAAAAGTTT	AAAAAGACCT	6180
ATCGTAGTCT	TGATCTTTTG	TTAATCGATG	ATAICCAATC	ACTCAGCGGA	AAAAAGTCG	6240
CAACTCAGGA	AGAAATTTTC	AATACCTTTA	ACGCCCTTCA	TGACAAGCAA	AAACAGATTG	6300
TCCTAACGAG	TGATCGTAGT	CCAAAACATC	TAGAAGGGCT	CGAGGAGAGG	CTGTCTACGC	6360
GTPTTAGTTG	GGGATGACA	CAAACTATCA	CCCCCCTGA	CTTTGAAACA	CGTATTGCCA	6420
TTTTACAAAG	TAAAGCGGAA	CAITTAGGCT	ACAATTTCCT	AAGTGATACT	CTAGAATACC	6480
TAGCTGGGCA	ATTTGATTTCA	AATGTTGAG	ATCTTGAGGG	AGCCATCAAC	GACATCACTT	6540
TAAATTGCCAG	AGTAAAAAAA	ATCAAGGATA	TCATATTGTA	TATTGCTGCA	GAAGCCATTA	6600
GAGCCCGCAA	ACAAGATGTT	AGCCAAATGC	TCGTATCCC	AATTGATAAA	ATCCAAACTG	6660
AAGTTGGTAA	CTTTTATGGT	GTTAGTATCA	AAGAAAATGAA	GGGAAGTAGA	CGCTTCAAA	6720
ATATTGTTTT	GGCCCGTCAA	GTAGCCATGT	ATTTATCTAG	AGAATAACA	GATAATATGC	6780
TTCCAAAAAT	TGGGAAGGAA	TTTGGGGGAA	AAGATCATAC	CACAGTCAIT	CATGCCCATG	6840
CCAAAATAAA	ATCTTTGATT	GATCAAGACG	ATAATTTACG	TTTAGAAAAT	GAATCAATCA	6900
AAAAGAAAAT	CAAAATAATT	GTGGATAACT	TTTAGTTTTT	TATCTTTTTT	ATCCACATTT	6960
TTTAAACAAG	CTAAAAAATC	TGATATGACT	TGTTTAAAGG	CTGTTTTCCA	CAGATTTTAC	7020
AGACTCTATT	ATTACTATTA	TCTTTCTAAT	ACTAAAAATA	AATAAGGAG	AATCCATGAT	7080
TCATTTTCCA	ATTAAATAAA	ATTATTTCT	ACAAGCATTA	AATACTACTA	AGAGAGCTAT	7140
TAGTTCTAAA	AATGCCATTC	CTATTTTATC	AACAGTAAAA	ATTGACGTGA	CCAATGAAGG	7200
TATTACTTTA	ATTGGTTCAA	ATGGTCARAT	TTCAATTGAA	AATTTTATTT	CTCAAAAAAA	7260
TGAAGATGCT	GGTTGTGTTA	TACTTCTTT	AGGTCGATCT	CTTCTTGAAG	CTTCTTTCTT	7320
TATCAATGTA	GTATCTAGTT	TACCTGATGT	AACCTTTGAT	TTTAAAGAAA	TTGAACAAAA	7380
TCAAAATGTT	TTAACCAAGT	GCAAAATCAGA	AATTAACCTTA	AAAGGAAAAG	ATAGCGAACA	7440
ATATCCACGA	ATCCAAGAAA	TTTCAGCAAG	CACCTCTTTA	ATACTTGAAA	CAAAATTACT	7500

262

CAAGAAAAAT	ATTAAATGAAA	CAGCCTTTGC	TGCAGTACACA	CAAGAGAGTC	GTCCGATTTT	7560
AACAGGGTTC	CACCTTCGTAT	TGAGTCAACA	CAAAGAGTTA	AAAACAGTTG	CAACAGACTC	7620
TCATCGCCTA	AGCCAGAAAA	AATTGACTCT	TGAAAAAAT	AGTGATGAT	TTGATGTCTG	7680
AATTCTTAGC	GCTTCTCTAC	GCGAATTTTC	AGCGGTATTT	ACAGATGATA	TCGAAACTGT	7740
AGAGATTTTC	TTTGCCAATA	ACCAAAATCCT	CTTTAGAAGC	GAAAAATATA	GCTTCTATAC	7800
TGCTCTCCTA	GAAGGAAACT	ATCCTGATAC	AGATCGCTTG	ATTCCAACAG	ACTTTAAACAC	7860
TACTATTACT	CTTAATGTGG	TAAACTTAGC	CCAGTCAATG	GAGCGTGCCC	GTCTTTTATC	7920
AAGTGCAGCT	CAAAATCGTA	CTGTGAAACT	TGAATTAAG	GATGGGGTTG	TTAGCGCCCA	7980
TGTTCACTCT	CCAGAACTTG	GTAAAGTAAA	CGAAGAAATC	GATACTGATC	AGGTTACTGG	8040
TGAAGATTTG	ACCATTAATT	TCAACCCAAC	TTACTTGATT	GATTCTCTTA	AAGCTTTTAA	8100
TAGCGAAAAG	GTGACTATTA	GCTTTATCTC	AGCTTGCTGT	CCATTTACTC	TTGTGCCAGC	8160
AGATACTGAC	GAAGACTTCA	TGCAGCTCAT	TACACCAGTT	COTACAAAT	AAGTGAAGA	8220
GOTTGAGCCT	GGCTCGCCTC	TTTTATGATA	TAATCGAAAA	AGAAAAGGAG	AGTAGTATGT	8280
ATCAAGTTGG	AAATTTTGT	GAGATGAAAA	AATCACACGC	TTGTACAATC	AAGTCGACTG	8340
GTAAAAAGGC	TAATCCTTGG	GAAATTACAC	GTGTAGGAGC	AGATATCAAA	ATAAAAATGTA	8400
GTAAATGTGA	GCATGTCTGC	ATGATGGGGC	GATATGATT	TGACCGAAAA	ATGAATAAAA	8460
TTATTGACTG	AGAACCCCTTA	GTAGAGGGT	TAGCACTTTA	TCCCTTTTTG	TGTTATAATA	8520
TTAGGGATTG	AAATGAAAAAC	GGAGAATGAG	AAATATGGCT	TTGACAGCAG	GTATCCTTGG	8580
TTTGCCAAC	CTTGCTAAAT	CAACACTATT	TATGCAATT	ACAAAAGCAG	GAGCAGAGGC	8640
AGCAAACTAC	CCATTTCCGA	CGATTGATCC	AAATGTTGGA	ATGCTGGAAG	TTCCAGATGA	8700
ACGCCTACAA	AAACTAACTG	AAATGATAAC	TCTTAAAAAG	ACAGTTCCCA	CAACATTTGA	8760
ATTTACAGAT	ATTGCAAGGA	TTGTAAAAGG	AGCTTCAAAA	GGAGAGGGGC	TAGGGAATAA	8820
ATCTCTGGCC	AATATTCTGT	AAGTAGATGC	GATTGTTTAC	GTAGTTCTGG	CTTTTGATGA	8880
TGAAATGTA	ATGCCCGAGC	AAGGACGTGA	AGACGCTTT	GTAGATCCAC	TTGCAGATAT	8940
TGATACCATT	AATCTGGAAT	TGATTCTTGC	TGACTTAGAA	TCAGTGAACA	AACGATATGC	9000
GCTGTAGAA	AAGATGSCAC	GTACCCAAAA	AGATAAGAA	TCAGTAGCAG	AACTCAATGT	9060
TCTTCAAAAG	ATTAAACCA	TCCTAGAAGA	CGGAAATCA	GCTCGTACCA	TTGAATTTAC	9120
AGATGAGGAA	CAAAAGGTTG	TCAAAGTCT	TTTCTTTTTC	ACGACTAAAC	CAGTTCTTTA	9180
TGTAGCTAAT	GTGGACGAGG	ATGTGGTTTC	AGAACCTGAC	TCTATCGACT	ATGTCAAACA	9240
AAATCTGTGA	TTTGCAGCGA	CAGAAAATGC	TGAAGTAGTC	GTATTTCTCT	CGCTGCTGA	9300

GGAGAAATTT	TCTGAATTGA	ATGATGAAGA	TAAAAAGAG	TTTCTTGAAG	CCATTGGTTT	9360
GACAGATATCA	GGTGTAGATA	AGTTGACCGG	TGCAGCTTAC	CACTTGGCTTG	GATTGGGAAC	9420
TTACTTCACA	GCTGGTGAAA	AAGAAGTTCC	CGCTTGGACT	TTCAACCTG	GTATGAAGGC	9480
TCTCAAGCA	GCTGGTATTA	TCCACTCAGA	CTTTGAAAAA	GGCTTTTATC	GTGCAGTAAC	9540
CATGTCATAT	GAAGACTAG	TGAAATAOAG	ATCTGAAAG	GCCGTAAAAA	AAGCTGGAGC	9600
CTTGGCTGAA	GAAGGAAAAA	AATATATCOT	TCAAGATGGC	GATATCATGG	AATTCCGCTT	9660
TAATGTCATA	AAATTAATAA	ATGGTGTCAA	TTAGGTTTGA	AAAAAATTC	AACCTTTTTC	9720
GCTTTTGAAA	GGAAAAATTA	ATGACCAAT	TACTTGTAGG	CTTGGGAAAT	CCAGGGGATA	9780
AAATATTTGA	AACAAAAAC	AATGTTGGTT	TTATGTTGAT	TGATCAACTA	GCGAAGAAAC	9840
AGAATGTCAC	TTTTACACAC	GATAAGATAT	TTCAAGCTGA	CCTAGCATCC	TTTTTCTCAA	9900
ATGGAGAAAA	AATTTATCTG	GTTAAACCAA	CGACCTTTAT	GAATGAAAGT	GGAAAAGCAG	9960
TTCATGCTTT	ATTAACTTAC	TATGGTTTGG	ATATTGACGA	TTTACTTTATC	ATTTACGATG	10020
ATCTTGACAT	GGAAGTGGG	AAAATTGCTT	TAAAGACAAA	AGGCTCAGCA	GGTGTCTATA	10080
ATGGTATCAA	GTCTATTATT	CAACATATAG	GAACTCAGGT	CTTTAACTGT	GTTAAGATTG	10140
GAAATTGAAG	ACCTAAAAAT	GGTATGTCAG	TTGTTTCACTA	TGTTTTCAAGT	AAGTTTGACA	10200
GGGATGATTA	TATCGGTATT	TTACAGTCTG	TTGACAAAGT	TGACGATTTCT	GTAAACTACT	10260
ATTACACAAG	GAAAAATTTT	GAGAAAAACA	TGCAGAGGTA	TAAACGATTA	ATGGTGACCT	10320
TATTAGATT	ATTCTCAGAA	AATGATCAGA	TTAAAAAATG	GCATCAAAAT	TTAACAGATA	10380
AGAAAAAGCA	ACTAATACTT	GGTTTATCAA	CATCTACTAA	GGCTCTTGCA	ATTGCAAGCA	10440
GTTTAGAAAA	AGAAGATAGG	ATTGTGTTAT	TGACGTCAAC	TTATGGAGAA	GCAGAAGGAC	10500
TTGTTAGTGA	TCTTATTCTT	ATCTTGGGTG	AGGAACCTGT	CTATCCATTT	TTGGTAGATG	10560
ATGCTCCTAT	GGTGGAGTTT	TTGATGCTTT	CACAGGAAAA	AATTATTTC	CGGGTTGAAG	10620
CCTTGGCTTT	TTTGAGTAT	TCACTCAAGA	AAGGGATT	AGTTTGTAAAT	ATCGCAGCAA	10680
GTCGATTGAT	TTTACCGTCT	CCCAATGCAT	TCAAAGATAG	TATTGTATAA	ATCTCAGTTG	10740
GTGAAGAACT	TGATCAACAC	GGCTTTATCC	ATCAGTTAAA	GGAAAAATGC	TATCGAAAG	10800
TTACTCAAGT	ACAAACTCAG	GGCGAATTTA	GTCTTCGAGG	AGATATTTTA	GATATTTTTC	10860
AAATATCCCA	GTTAGAACCT	TGTCGAATTG	AGTTTCTTGG	TGATGAAATT	GATGGTATCA	10920
GGTCATTGTA	AGTAGAACA	CAATTAATCGA	AAGAAAAATA	GACAGAACTC	ACTATCTTTC	10980
CAGCTAGTGA	TATGCTTTTG	AGAGAAAAAG	ATTATCAACG	AGGACGATCA	GCTTTAGAAA	11040

264

AACAAATTC AAAACTTTA TCACCTATTT TGAATTCATA CCTAGAAGAA ATTCCTITCAA	11100
GTPTTCACCA AAAACAAGT CATGCAGACT CTCGGAAGTT TTTATCTTTG TGTATGATA	11160
AGACATGGAC TGCTTTGAT TATATTGAAA AAGATACTCC AATATTTCTT GATGATTATC	11220
AAAAACTGAT GAATCAGTAT GAAGTCTTTG AAAGAGACTT AGCGCAGTAC TTACAGAAG	11280
AATTACAGAA TAGTAAAGCA TTTTCTGATA TGCAGTATTT TTCTGATATT GAACAAATCT	11340
ATAAAAAACA AAGTCCAGTG ACCTTTCTCT CTAATCTTCA AAAGGGTTTA GGAATCTCA	11400
AATTGACAAA AATTATCAAA TTCAATCAAT ATCCTATGCA GGAATTTTTC AATCAGTTTT	11460
CTTTCTTAAA AGAAGAAAT GAACGATATA AAAAATGGA TTACACCATT ATTCCTGAGT	11520
CTAGCAATTC AATGGGAAGT AAAACATGAG AGGATATGTT AGAGGAATAT CAGATTAAAT	11580
TGGATCTAG AGATAAGACA AATATCTGTA ANGAATCTGT AACTTAATA GAGGGTAATC	11640
TCAGACATGG TTTTCATITT GTAGATGAAA AGATTTTATT GATAACTGAA CATGAGATTT	11700
TTCAAAGAAA ATTAAAGGT CTTTTCGAA GACAACATGT TTCAATGCA GAGAGATTAA	11760
AAGATTACAA TGAACCTGAA AAAGGGGACT ATGTTGTCCA TCATATCCAT GGGATTGGTC	11820
AATATCTAGG AATTGAAACC ATTGAAATCA AGGGAATTCA TCGGATTAT GTCAGTGCC	11880
AATACCAAAA TGGTGATCAA ATTTCTATCC CCGTGAACA GATTATCTA CTGTCCAAAT	11940
ATATTTCAAG TGATGGTAAA GCTCCAAAAC TCATATAAT AAATGACGGT CATTTTAAAA	12000
AGGCCAAGCA AAAGGTTAAG AACCAGGTAG AGGATATAG TGATGATTAA ATCAAACTCT	12060
ACTCTGAACG TAGTCAGTTG AAGGTTTTG CTTTCTCAGC TGATGATGAT GATCAAGATG	12120
CCTTTGATGA TGCTTCCCT TATGTTGAAA CGGATGATCA ACTTCGTAGT ATTGAGGAAA	12180
TCAGAGGGA ATGCAAGGCT TCTCAGCCAA TGGATCGACT TTTAGTTGGG GATGTTGGTT	12240
TTGGAAGAC TGAAGTTGCT ATGCGTGCAG CCTTTAAGC AGTCAATGAT CACAAACAGG	12300
TTGTCAATCT AGTCCCGACG ACGGTTTTAG CGCAACAGCA CTATACGANT TTTAAGGAAC	12360
GATTCACAAA TTTTGCAGTT AATATTGATG TGTGAGTCG CTTTAGAAGT AAAAAGAGC	12420
AGACTGCAAC ACTTGAAAAA TTGAAAAAGC GTCAAGTCGA TATTTTGATT GGAACACATC	12480
GTGTTTTGTC AAAAGATGTT GTGTTTGCTG ATTTGGGCTT GATGATTATT GATGAGGAAC	12540
AGCGATTTCG TGTCAGCAT AAGGAAACTT TGAAGAACT GAAGAAACAA GTGGATGTC	12600
TAACCTTGAC CGTCAAGCCA ATCCCTCGTA CCTTCCATAT GTCTATGCTG GGAATCAGAG	12660
ATTTATCTGT TATTGAAACT CCGCGACTA ATCGCTATCC TGTTCAGACC TATGTTTTGG	12720
AAAAGAATGA TAGTGTCATT COTGATGCTG TCTTGCGTGA AATGGAGCGT GGAGGTCAAG	12780
TTTATATCT TTACACAAA GTTGACACAA TTGTTCAGAA GGTTCAGAA TTACAGGAGT	12840

TGATTCGGA	GCCTTCGATT	GGATATGTTT	ATGGTCGAAT	GAGTGAAGTC	CAGTTGGAAA	12900
ATACTCTATT	AGACTTTATT	GAGGACAAAT	ACGATATCTT	GGTGACGACT	ACTATTATTG	12960
AGACAGGGGT	GGACATTCCA	AATGCTAATA	CTTTATTAT	TGAAAATCG	GACCATATGG	13020
GCTTGTCAAC	CTTATATCAG	TTAAGAGGAA	GAGTCGTCG	TAGTAATCOT	ATTGCTTATG	13080
CTTATCTCAT	GTATCGTCCA	GAAAAATCAA	TCAGTGAAGT	CTCTGAAAAG	AGATTAGAAAG	13140
CGATTAAAGG	ATTTACAGAA	TTGGGCTCTG	GCTTTAAGAT	TGCAATGCCA	GATCTTTGCA	13200
TTGCTGGAGC	AGGAAATCTT	TTAGGAAAAAT	CCGAGTCTGG	TTTCATTGAT	TCGTGTTGGT	13260
TTGAAATTGA	TTGCGAGTTA	TTAGAGGAAG	CTATTGCTAA	ACGAAACGGT	AATGCTAAGC	13320
CTAACACAAAG	AACCAAAGGG	AATGCTGAGT	TGATTTTGCA	AATTGATGCC	TATCTTCTCG	13380
ATACTTATAT	TTCTGATCAA	CGACATAAGA	TTGAAATTTA	CAAGAAAAAT	CGTCAAAATTG	13440
ACAACCGTGT	CAATTATGAA	GAGTTACAAAG	AGGAGTTGAT	AGACCGTTTT	GGAGAAATACC	13500
CAGATGTAGT	AGCCTATCTG	TTAGAGATTG	GTTTGGTCAA	ATCATACTTG	GACAAGGTCT	13560
TTGTTCAACG	TGTGGAAGA	AAAGATAAAT	AAATTACAAT	TCAATTTGAA	AAAGTCACTC	13620
AACGACTGTT	TTTAGCTCAA	GATTATTTTA	AAGCTTTATC	CGTAACGAAC	TTAAAAGCAG	13680
GCATCGCTGA	GAATAAGGA	TTAATGGAGC	TTGTATTTTGA	TGTCCAAAAT	AAGAAAGATT	13740
ATGAAATTTT	AGAAGGTTTG	CTGATTTTTC	GAGAAAGTTT	ATTAGAGATA	AAAGAGTCTA	13800
AGGAAGAAAA	TTCCATTGTA	TATTTTTCTT	CTATAAAATA	GATRAAAATG	GTACAATAAT	13860
AAATTGAGGT	AATAAGGATG	AGATTAGATA	AATATTAAAA	AGTATCCGGA	ATTATCAAGC	13920
GTGCTACAGT	CGCAAAGGAA	GTAGCAGATA	AAGGTAGAAT	CAAGGTTAAT	GGAATCTTGG	13980
CCAAAGGTTT	AACGAGCTTG	AAAGTTAATG	ACCAAGTTGA	AATTGCTTTT	GGCAATAAGT	14040
TGCTGCTTGT	AAAAGTACTA	GAGATGAAGG	ATAGTACAAA	AAAAGAAGAT	GCAGCAGGAA	14100
TGTATGAAAT	TATCAGTGAA	ACACGGGTAG	AAGAAAAATG	CTAAAAATAT	TGTACAATTG	14160
AAATAATTCTT	TTATTCAAAA	TGAATACCAA	CGTCGTGCTT	ACCTGATGAA	AGAACGACAA	14220
AAACGGAATC	GTTTTATGGG	AGGGGTATTG	ATTTTGATTA	TGCTATTAT	TATCTTGCCA	14280
ACTTTTAATT	TAGCGCAGAG	TTATCAGCAA	TTACTCCAAA	GACGTGACGA	ATTAGACAGC	14340
TTGCARATC	AGTATCAAA	TTTGAAGTAT	GAAAGGATA	AGAGACAGC	ATTGCTTACC	14400
AAGTTGAAG	ATGAAGATT	TGCTGCTAAG	TATACAGAG	CGAAGTACTA	TTATTCTAAG	14460
TCGAGGGAAA	AAGTTTATAC	GATTCCTGAC	TTGCTTCAAA	GGTGATAAAA	TGAAAAATTT	14520
ATTAGACGTA	ATAGAGCAAT	TTTTGAGTTT	GTGAGATGAA	AAGCTGGAAG	AATTGGCTGA	14580

256

TAAAAATCAA	TTAATGCGTT	TACAAGAAGA	AAAGGAAGG	AAGATGCCGT	AAATCTTTAA	14640
TTAATTTGTT	GCTACCAAGT	TTTTTGACCA	TTTCAAAAGT	CGTTAGCACA	GAAAAAGAAG	14700
TCGTCTATAC	TTGGAAGAAA	ATTATTATACC	TTTCACAAATC	TGACTTTGGT	ATTTATTTTA	14760
GAGAAAAATT	AAGTCTCCCT	ATGGTTTATG	GAGAGGTTCC	TGTTTATCGG	AATGAAGATT	14820
TAGTAGTGGG	ACTGGGAGAA	TTGACTCCCA	AAACAAGTTT	TCAAATAACC	GAGTGGCGCT	14880
TAAATAAACCA	AGGAATTCCT	GTATTTAAGC	TATCAAAATCA	TCAATTTATA	GCTGCGGACA	14940
AACGATTTTT	ATATGATCAA	TCAGAGGTAA	CTCCAACAAT	AAAAAAAGTA	TGGTTAGAAT	15000
CTGACTTTAA	ACTGTACAAT	AGTCTTTATG	ATTTAAAAGA	AGTGAAATCA	TCCTTATCAG	15060
CTTATTTCGA	AGTATCAATC	GACAAGACCA	TGTTGTGAGA	AGGAAGAGAA	TTTCTACATA	15120
TTGATCAGCG	TGGATGGGTA	GCTAAAGAAAT	CAACTTCTGA	AGAAGATAAT	CGGATGAGTA	15180
AAGTTCAGGA	AATGTTATCT	GAAAAATATC	AGAAAGATTC	TTTCTCTATT	TATGTTAAGC	15240
AAC TGACTAC	TGGAAAGAGAA	GCTGGTATCA	ATCAAGATGA	AAAGATGTAT	CGAGCCAGCG	15300
TTTTGAAACT	CTCTTATCTC	TATTATACGC	AAGAAAAAAT	AAATGAGGGT	CTTTATCAGT	15360
TAGATACGAC	TGTAAAATAC	GTATCTGCCG	TCAATGATTT	TCCAGGTTCT	TATAAACCCAG	15420
AGGGAAGTGG	TAGTCTTCTT	AAAAAAGAAG	ATAATAAAGA	ATATTCTTTA	AAGGATTTAA	15480
TTACGAAAGT	ATCAAAGAGAA	TCTGATAATG	TAGCTCATAA	TCTATTGGGA	TATACATTTT	15540
CAAAACCAATC	TGATGCCACA	TTCAAATCCA	AGATGTCTGC	CATTATTGGGA	GATGATTGGG	15600
ATCCAAAGAA	AAATTTGATT	TCTTCTAAGA	TGCGCCGGAA	GTTTATGGAA	GCTATTTTATA	15660
ATCAAAATGG	ATTTGTGCTA	GAGTCTTTGA	CTAAAACAGA	TTTTGATAGT	CAGCGAATFG	15720
CCAAAGGTGT	TTCTGTAAAA	GTAGCTCATTA	AAATTGGAGA	TGCGGATGAA	TTTAAGCATG	15780
ATACGGGTGT	TGCTATGACA	GATTCTCCAT	TTATTCTTTC	TATTTTCACT	AAGAATCTGT	15840
ATTATGATAC	GATTTCTAAG	ATAGCCAAGG	ATGTTTATGA	GTTTCTAAAA	TGAGGGAACC	15900
AGATTTTTTA	AATCATTTTC	TCAAGAAGGG	ATATTTCAAA	AAGCATGCTA	AGGCGTTTCT	15960
AGCTCTTTCT	GTTGGATTAG	ATTCCAATGT	TCTATTTAAG	GTAATGTCTA	CTTATCAAAA	16020
AGAGTTAGAG	ATTGAATTGA	TTCTAGCTCA	TGTGAATCAT	AAGCAGAGAA	TTGAATCAGA	16080
TTGGGAAGAA	AAGGAATTAA	GGAAGTTGGC	TGCTGAAGCA	GAGCTTCTTA	TTTATATFCA	16140
CAATTTTTCCT	GGAGAATTTT	CAGAAGCGCG	TGCACGAAT	TTTTGTTATG	ATTTTTTTCA	16200
AGAGGTCATG	AAAAAGACAG	GTGCGACAGC	TTTAGTCACT	GCCCACCATG	CTGATGATCA	16260
GGTGAAACG	ATTTTATATG	GCTTGATTCG	AGGAACTCGC	TTGCGCTATC	TATCAGGAAT	16320
TAAGGAGAAG	CAAGTAGTCG	GAGAGATAGA	AATCATTCGT	CCCTTCTTGC	ATTTTCAGAA	16380



AAAAGACTTT CCATCAATTT TTCACITTTGA AGATACATCA AATCAGGAGA ATCAITATPTT 16440  
 TCGAAATCGT ATTGCGAAATTT CTACTTACC AGAATTGGAA AAAGAAATC CTCGATTTAG 16500  
 GGATGCAATC TTAGGCATTG GCAATGAAAT TTTAGATTAT GATTTCGGCA TAGCTGAAT 16560  
 ATCTAACAT ATTAATGTGG AAGATTTTACA GCAGTTATTT TCTTACTCTG AGTCTACACA 16620  
 AAGAGTTTTA CTTCAAACTT ATCTGAATCG TTTTCAGAT TTGAATCTTA CAAAGCTCA 16680  
 GTTTGCTGAA GTTCAGCAGA TTTTAAATC TAAAGCCAG TATCCTCATC CGATTAAAA 16740  
 TGGCTATGAA TTGATAAAAG AGTACCAACA GTTTCAGATT TGTAAATCA GTCCGCAAGC 16800  
 TGATGAAAG GAAGATGAAC TTGTGTTTACA CTATCAAAAT CAGGTAGCTT ATCAAGGATA 16860  
 TTTATTTTCT TTTGCACTTC CATTAGAAGG TGAATTAAAT CAACAAATAC CTGTTTCACG 16920  
 TGAACATCC ATACACATTC GTCATCGARA AACAGGAGAT GTTTTGATTA AAAATGGGCA 16980  
 TAGAAAAA CTCAGACGTT TATTATTGTA TTGAAAACT CCTATGGAAA AGAGAAATCT 17040  
 TGCTCTTATT ATTGAGCAAT TTGTGAAAT TGTCTCAAT TTGGGAATPG CGACCAATAA 17100  
 TTTGATTAAT AAAACGAAAA ATGATATAAT GAACACTGTA CTTTATATAG AAAAAATAGA 17160  
 TAGGTAAAA ATGTTAGAAA ACGATATTAA AAAAGTCTCT GTTTCACAG ATGAATTTAC 17220  
 AGAAGCAGCT AAAAACTAG GTGCTCAATT AACTAAGAC TATGCAAGAA AAAATCCAAT 17280  
 CTTAGTTGGG ATTTTAAAG GATCTATTCC TTTTATGCT GAATTGGTCA AACATATTGA 17340  
 TACACATATT GAAATGSACT TCATGATGGT TTCTAGCTAC CATGGTGGAA CAGCAAGTAG 17400  
 TGGTGTATAC AATATTAAAC AAGATGTGAC TCAAGATATC AAAGGAAGAC ATGTTCTATT 17460  
 TGTAGAAGAT ATCAATTGATA CAGTCAAAAC TTTGAAGAA TTGCGAGATA TGTTTAAAGA 17520  
 AAGAGAAGCA GCTTCTGTTA AAATGCAAC CTGTTGGAT AAACGAGAAG GACGTGTTGT 17580  
 AGAAATTGAG GCAGACTATA CTGCTTTTAC TATCCCAAT GAGTTGTAG TAGGTTATGG 17640  
 TTTAGACTAC AAAGAAATTT ATCGTAATCT TCCTTATATT GGAGTATTA AAGAGGAAGT 17700  
 GTATTCAAAAT TAGAAGAAAT AATCTTTTAAAT GAAAAAACAA AATAATGGTT TAAATAAAA 17760  
 TCCITTTCTA TGGTTATTAT TTATCTTTT COTGTGACA GGATTCCAGT ATTTCTATTC 17820  
 TGGGAATAC TCAGGAGGAA GTCAGCAAT CAACATATCT GAGTTGCTAC AAGAAATTAC 17880  
 CGATGGTAAT GTAAAGAAAT TAACTTACCA ACCAATGGT AGTGTATTAG AAGTTTCTGG 17940  
 TGCTATAAAA AATCTTAAAA CAAGTAAAGA AGAAACAGGT ATTCAGTTT TCACGCCATC 18000  
 TGTTACTAAG GTAGAGAAAT TTACCAGCAC TATCTTCTCT GCAGATACTA CCGTNTCAGA 18060  
 ATTGCAAAAA CTGTCTACTG ACCATAAAGC AGAAGTAAT GTTAAAGCATG AAGTTCAAG 18120

268

TGCTATATGG ATTAATCTAC TCGTATCCAT TGTGCCATTT	GGAATCTAT TCTTCTCTCT	18180
ATPCTCTATG ATGGGAAATA TGGGAGGAGG CAATGGCCOT	AATCCAATGA GTTTTGGAGC	18240
TAGTAAGGCT AAAGCAGCAA ATAAAGAAGA TATTAAGTA	AGATTTTCAG ATGTTCCTGG	18300
AGCTGAGGAA GAAAAACIAG AACTAGTTGA AGTTGTGTAG	TTCTTAAAG ATCCAAAACG	18360
ATTCACAAAA CTGGAGCCC GTATTCACAG AGGTGTCTT	TTGGAGGAC CTCGGGGAC	18420
AGGTAAAJCT TTGCTTGCTA AGGCAGTCCG TGGAGAAGCA	GGTGTTCAT TCTTTAGTAT	18480
CTCAGGTCT GACTTTGTAG AAATGTTTGT CGGAGTTGGA	GCTAGTCGTG TTCGCTCTCT	18540
TTTTGAGGAT GCCAAAAAG CAGCACGAG TATCATCTTT	ATCGATGAAA TTGAATGCTGT	18600
TGGACGTCAA CGTGGAGTCG GTCTCGGCGG AGGTAAATGAC	GAACGTGAAC AAACCTTGAA	18660
CCAACTTTTG ATTGAGATGG ATGGTTTTGA GGGAAATGAA	GGGATTATCG TCATCGCTGC	18720
GACAAACCGT TCAGATGTAC TTGACCTTGC CCTTTTGGGT	CCAGGACGTT TTGATAGAAA	18780
AGTATGGTT GGTGCTCCTG ATGTTAAAG TCGTGAAGCA	ATCTTGAAG TTCACGTAA	18840
GAATAAGCCT TTAGCAGAAG ATGTTGATTT GAAATTAGTG	GCTCAACAAA CTCAGGCTT	18900
TGTTGCTGCT GATTTAGAGA ATGCTTTGAA TGAAGCAGT	TTAGTTGCTG CTGCTGCCTA	18960
TAAATCGATA ATTGATGCTT CAGATATTGA TGAAGCAGAA	GATAGAGTTA TTGCTGGACC	19020
TTCTAAGAAA GATAAGACAG TTTCACAAAA AGAACGAGAA	TTGTTGCTT ACCATGAGGC	19080
AGGACATACC ATTGTTGCTC TAGTCTTGTC GAATGCTCGC	GTGTTCCATA AGGTTACAAT	19140
TGTACCACGC GGCCGTGCAG GCGGATACAT GATTGCACTT	CCTAAAGAGG ATCAAATGCT	19200
TCTATCTAAA GAAGATATGA AAGAGCAATT GCGTGCCTTA	ATGGGTGGAC GTGTAGCTGA	19260
AGAAATTATC TTTAATGTCC AAACCACAGG AGCTTCAAAC	GACTTTGAAC AAGCGACACA	19320
AATGGCACGT GCAATGTTA CAGAGTACGG TATGAGTGAA	AAACTTGGCC CAGTACAATA	19380
TGAAGGAAC CATGCTATGC TTGGTGCACA GAGTCCCTCA	AAATCAATTT CAGAACAAAC	19440
AGCTTATGAA ATTGATGAAG AGGTTGCTTC ATTATTAAAT	GAGGCACGAA ATAAAGCTGC	19500
TGAATTATTT CAGTCAAAATC GTGAAACTCA CAAGTTAAT	GCAGAAGCAT TATTGAJATA	19560
CGAAJACTTG GATGATACAC AAATTAAGC TCTTTACGAA	ACAGGAAGA TGCCGTGAAGC	19620
AGTAGAAGAG GAATCTCATG CACTATCCTA TGATGAAGTA	AAGTCAAAAA TGAATGACGA	19680
AAATAAACCC TGAGAGAGGC TGGAGCCTCT CTTTTTGTG	GCTTATAGGA GCTAAAGGGA	19740
ACAGAATGGA GAAAAAGGAA CAAATGTGTT TTCTAATCTG	TTAGACTGTA TCTAGAAAGG	19800
GGAAATTTAT GATTAAGAA TTGTATGAAG AAGTCCAAGG	GACTGTGTAT AAGGTAGAA	19860
ATGAATATTA CTTTCAATTA TGGGAATTGT CGGATTGGGA	GCAAGAAGGC ATGCTCTGCT	19920

TACATGAATT GATTAGTAGA GAAGAAGCAC TGTAGACGA TATTCCACGT TTAAGGAAT 19980  
 ATTTCAGAC CAAGTTTGA AATCGAATT TAGACTATAT CCGTAACAG GAAAGTCAGA 20040  
 AGCOTAGATA CGATAAAGAA CCCTATGAAG AAGTGGGTGA GATCAGTCAT COTATAAGTG 20100  
 AGGGGGGTCT CTGGCTAGAT GATTATTATC TCTTCATGA AACACTAAGA GATTATAGAA 20160  
 ACAACCAAG TAAAGAGAA CAAGAAGAAC TAGAACCGT CTTAAGCAAT GAACGATTTC 20220  
 GAGGGCGTCA AAGAGTATTA AGAGACTTAC GCATTGTGTT TAAGGAGTTT ACTATCCGTA 20280  
 CCCACTAGTA AGTCATGCAA AAAAAATGAA AAAAATTAGA AAAAGTAGTT GACAAAGTTT 20340  
 GAAAAGGCTG TATAATAGTA AGAGTTGAAA ATAAACACTC AGTCCGTTG GTCAGGGGT 20400  
 TAAGACACCG CCTTTTCAG CCGTAACAC GCGTTCGAAT CCCGTACGGA CTATGGTATG 20460  
 TTGCGTCAGG ACCACTTGAT GAAAAAAGT TTAATAAATCT TCAAAAAGT 20520  
 GTTGACAAGC GAAAGCAGTT GTGATATACT AATATAGTTG TCGCTTGAGA GAAGCAAGTG 20580  
 ACAAGACCT TTGAAAACCT AACAGACGA ACCAATGTGC AGGCGCTAC AACGTAAAGTT 20640  
 GTAGTACTGA ACAATGAAAA AAACAATAAA TCTCTCATGT ACAGAAATGA GTAAGAACTC 20700  
 AAATTTTTTA ATGAGAGTTT GATCCTGGCT CAGGACGAAC GCTGGCGCGG TGCCTAATAC 20760  
 ATGCAAGTAG AACGCTCAAG GAGGAGCTTG CTCTCTGGA TGAGTTGCGA ACGGGTGAGT 20820  
 AACGCGTAGG TAACCTGCCT GTAGCGGGG GATAACTATT GGAAAGATA GCTAATACCG 20880  
 CATAAGAGTA GATGTTGCAT GACATTGCT TAAAGGTGC ACTTGCAATCA CTACAGATG 20940  
 GACCTGCTT GTATTAGCTA GTTGCTGGG TAACGGCTCA CCAAGCGAC GATACATAGC 21000  
 CGACCTGAGA CGGTGATCG CCACACTGGG ACTGAGACAC GCGCCAGACT CCTACGGGAG 21060  
 GCAGCAGTAG GCAATCTTCG GCAATGGACG GAAGTCTGAC CGAGCAACGC CGCGTGAGTG 21120  
 AAGAAGGTTT TCGAGCTGTA AAGCTCTGTT GTAAGAGAAG AACGAGTGTG AGAGTGGAAA 21180  
 GTTCACACTG TGACGATATC TTACCAGAAA GGGACGCTA ACTACGTGCC AGCAGCGCGG 21240  
 GTAAATACATA GGTCCCGAGC GTTGTCGGGA TTTATTGGCG GTAAAGCGAG CGCAGCGGT 21300  
 TAGATAAGTC TGAAGTTAAA GGCTGTGGCT TAACCATA 21338

(2) INFORMATION FOR SEQ ID NO: 21:

(i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 6273 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

270

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 21:

TGTTTTTAAAGAGCCGTGTC	TGGATAGACT	TTCCGACGCA	ACGCTCTATT	AGATAATGAA	60
CTGCCATATAC	ACRAGATTTC	TAACCTTAGT	CGACATGAGC	TGAACCTCT	120
GTAGTTCACA	AAATATTATA	CACCTATTTT	ATGAMTAGTC	AACGTGCTTT	180
TTTGAAGAAAT	CATGAAAAT	TTCTCTTTCT	TTCCATTTTA	AGTGACATTC	240
ACATCAJAAA	AGCCCAGAGC	AAATTGCTG	AGCATCTTT	TATCTAGTCG	300
TTGAGTTCAG	TATGTTTAAA	GTCTCTGTCC	CATCATTTCT	TCAACAAACC	360
AGAAACTCCT	TGGCTACTTG	CTTTGCTGAC	TTGCCCTTCA	CACCGACTTG	420
TGGCTCATCT	GGCTTCTGT	AATCTTACCA	GCCAATGTAT	TAGAACTCT	480
GGGTGTTTCT	TGAGAAGAGC	TTCTTTTCA	AGTGGAGCCC	CTTGATAAGG	540
TGCTTGTTCAT	CTTCCAAGAC	CTGTAAATCA	TAAAGCTCCA	ATTCGGCATC	600
GCATCCGTGA	TTTGAATATC	CCCTGACTGA	ATAGCCGTAT	AGCGAAGGGC	660
GTCCGTACAT	TGAGATTGAG	ACCATACATT	GATTGCAAGC	CCTTATTTCC	720
TCGTAAACT	CGAGTGTAAA	ACCTGCCCTC	AACGCCCC	CCACTTTTTT	780
ATGCTCTTCA	AGCCATATTC	TTGAGCAATC	TTTTTCOGAA	CAGCTACAGC	840
TGATAGACA	TGGGTTTGAG	ATAGGCTAGA	TGATCTGTCT	TAGCAATGCC	900
ACCTGATAAA	CCTGTTCTGG	TTCAATGACTC	ACCTTGGGTG	ATGGTTGAAG	960
GTCACCGTAC	CAGTAAATTC	AGGATAGATG	TCAATATGCC	CTTTTTCAG	1020
AGGAAGCTTG	TCTTCCCAAA	ATTCGGTTTA	ACAGTCGCAG	TCATGCTGGT	1080
ATCAGCAACT	TATACATATT	GGCCAAAATT	TCTGGTCTGT	GACCTATTTT	1140
ACCAAGTTTT	CCTTCTCTTT	TTGAACCAAA	AGAGCTGGAC	TATAAGACAG	1200
AAAGCCACCA	AGGCAAAACC	TGAGAAAATC	GTCCGTAATT	TTGCTTTTTTC	1260
AGTAGGAAGT	TAAAGGCAAT	GGCTAGCACT	CGAGAGAGAA	GTGCCCAAT	1320
CTGGCATAT	TACGCTCAAT	TCCCAAAAGA	ATAAAGGAAC	CTAGTCCCCC	1380
AGGCCGCCCA	AGGTTGCGGT	ACCGATAATC	AAAAACGCTG	CCGTCCGAAT	1440
ATAACAGGCA	TGGCGAGTGG	AATTTCAAAT	TTCTTGAGAC	GTTCCTCATCT	1500
AAGGCAATCC	CAGCCTCTTG	CAGGTTGGGA	TCAATTCCCT	TCAGGCCAGT	1560
TGCAAAATAG	GGAAAATGCG	ATAAATCACT	AGAGCTGTCA	AAGCCCGCAA	1620
CCCATCAAG	GGATAAAGAG	CCCCAACAG	GCCAGAGAGC	GGATGCTCTG	1680
GCAATCTGCA	AGACCCAGTC	GGCCAGCTTC	TCATGATAGC	GAAGAAAAC	1740

ATCGCAAGCA	AAATAGCTAG	TAAACAAGGTC	AAAAGCGACA	ACTGCAATG	TTGAGATAGA	1800
GCTGTCAACC	AATCACTAAA	ACGATCCTGA	AAGTTGCAA	TTAAATTAGT	CATGAACACT	1860
ACCTCCAAAC	AAGTCTGCTA	CAAAGTCTGT	TGCAGGCGCT	TTTAAAAATTG	TCTCGGGATT	1920
CGCTACCTGG	CGAATTCTCT	CATCCTGCAA	GACAGCAATA	CGGTCCGCCA	ACTTCAAGGC	1980
TTCACTCGTA	TGATGGGTTA	CAAAAATCGT	TGTCATCCCA	AACCTCTTAT	GCAATCTCTT	2040
TGTCAGAAC	TGCAACTGTT	TTCTCGAANT	AGCATTCCAAG	GCCGAAAAGG	GTTCAATCCAT	2100
GAGGAAAATC	TTGGGCTGAC	CAATCATAGC	TCGGACAATA	CCGACCCGTT	GCTGTTCTCC	2160
ACCAGATAAT	TCACTAGGTA	AGCGATGCCC	ATACTCGGCT	ACTGGTAAC	CAACCTTAGC	2220
CAAAAGCTCT	TCTGTTTTCT	TCGTAATTC	TTCCCTTGCTC	CACCCCTTCA	TTTCAGGAAT	2280
GAGAGCAATA	TTTTCCGCAA	CTGTTAGATT	TGGAAAAGA	GCAATAGCCT	GTAAACATA	2340
ACCAGTAGAA	AGACGAAATT	CACGCTCATC	ATAGTCTTTG	ATCGGCTTCC	CATCCATATA	2400
AATATTTCOA	TCAGTTGGTT	CCAAAAGACG	GTAAATCATC	TTGAGCATGG	TGCTCTTACC	2460
TGACCCAGAA	GGCCCTACTA	AAACCATAA	TTCCCATCC	TCAATCTGTA	AGTTGACATC	2520
TCTCAAGACA	TCTTTTCTG	TGTAGCGCAG	TGCTACATTT	TTGTATTCAA	TCATTCTTTG	2580
TCCTCAATTT	AAAACCTCCC	TCGATTGGTC	AAGTCTCTTA	CCTTAGGCAT	AACCTTCCTTA	2640
TATATCCCAAT	GCTCCACAAT	TTTCCGTTTC	TCATAAACGA	AGATATCGTA	CTGGGCAATA	2700
GCAACGCCAT	CAATCTGAST	CTGACCATAG	CTAACCATAT	AGTTTCTCTG	TCCTAAGAGT	2760
TGGAAAACAA	AGTCAAAAGT	GACACTATAT	TCAGCCACAT	AGTTTTTATA	AGCAGCACTT	2820
CCTTGTCCAA	TATCATGATT	ATGCTGAATC	AAATCGTCTG	CCACATAATC	ACTCCACTGC	2880
TCTAGCTCCC	CATTTTGGAA	AATTTCTGTC	AAGAAACGGC	GAACCAGCTT	TTTATTTTCT	2940
GCTTTCCTAT	CCAAATCCTT	GATTTCAAAA	TCTCAAAAA	TTTGATCTAG	TTGGTCAATT	3000
TCAGGTGTTT	GATAGTAGTC	AATGACATCC	CAATGCTCAA	CAATACAAAC	ATTTCTCATCC	3060
TCACGGAAG	TATCCGTCGT	CACCCATTGA	GCTTCTCCAC	CATTGAGATA	TTGATGAACA	3120
TGAACAAAGA	CCAGATTGCC	ATCCTCAATG	GTGCGGACAA	TCTTAATCTG	ACGCTCTGSA	3180
TGACGCTCAA	AGAAATCTGC	AAAGAAGGCT	GCAATCCTT	CTTTCCCGTC	AGGAACACCT	3240
GTCAATGTT	GGATATAGGT	ATCCCTTACA	GACTGGGCTT	GAGCTCAGC	AACCTGTCOG	3300
TCTTGAATGG	CATGAGATGA	TAGGTTGTGA	GCATTTTICA	CTTGTGTGGA	CATATTCTAA	3360
ACCTCATTTT	CCTTCTCTTT	CAGATTGCGC	AAAAATCTTT	CTTGAAAC	TTCAANTTGG	3420
TGAATTTCTT	CCTCTGAAAA	TCCTTTGTAA	AAGATAGTAT	CCAAATTTCTG	ACTGACACGA	3480

272						
TGCCCACTT	CTPTTCGGGA	CTTGCCCTAAC	TCCGTTAAAA	CTAAATACTT	CTTACGCTTG	3540
TCTTTTCCAC	ACGGCATAAC	AATTACAAGC	TTTGTGTCTCT	CTAGCTTTTTT	TATCATAGTC	3600
GTCAAGCTAT	TATTTCGCAAG	TCCAGTTCGCA	AGCGCGATAT	CTGTGCCGAGT	TGCGCAGCCA	3660
GTTTCACTAT	TCCATAAAAC	CGCTAAAATC	TTGCCCTGTT	CACCCCTATA	AAGAGCCTCA	3720
GGATCTTGAC	TCAGTAACCT	TTGAAAAATC	CGCCCATATCA	ACAAACGAAT	ATGATGGGCT	3780
AGCAAAATGAC	CATCTTTCAT	AACACCTCCA	ATTTATTTCG	ATATCGAAAT	GAATAAAACA	3840
ATTGTAAACAC	TCATCGTTCT	AACTGTCAAC	TATTTGATG	TAGAAATAAT	TTTTGATAAT	3900
TATCCACACC	ACCACTCTCC	GGCTCAACTA	ACTTTTAAAG	AGAGTTTCTA	AACTCCTTCG	3960
TCCTCCAGTC	TACAAAAGCC	TTCCATTGCT	ACTATCCTAT	ATTTTATGAG	GGGACACATT	4020
TTTCTCATCA	GACCATTTAT	TTTAAAGATA	GAAGTAAATC	ATAATTGCTT	CCATCTGTTC	4080
TTTTATAGTA	TATTGAAGTT	AGACTAGAGC	ACTGTATCTT	CTAAACANTT	GATAGAAAGC	4140
GATTTGAATT	TCCCAATCAA	TTTGTTCGTA	TTTATAGCAT	TTCGAAACTG	GAATAGGACA	4200
CCATGACTGC	TAAAAGATTT	CTATAAATTC	ATTTAATTTC	CTCAATCAAT	TTGTTCATAT	4260
CTTATTTTCAT	TCCGCTATAA	TTTCACTTTA	CCCTATCTTT	TTTGTAGCAC	CCTTCAAACA	4320
GCCTATCCCC	TACCGTTTGA	CGATTTCTCA	CTTCGCTCCA	CTTCCATTAC	AGAAGTTTCT	4380
TCACTACTAT	GGGCTCGGCT	GACTTCTCAT	GATTCCCTGT	TACTACTATT	TGAACGCTCA	4440
CGAGATAGAT	CTTACAAAAA	ATGCTTTGAT	CCACAATGGA	ATCAAAGCAT	TTTAAAGAGT	4500
TCCTCATACA	TAAAGCGAGA	AGTCGCGAGT	CCTCTGTACT	TGGCTTCTTC	TCTTTTGACA	4560
AAGCGAGCCA	AGTTGAGCAA	CTCAGGTGCT	GGATGTTTGG	GATTTAGGAG	CAATTCACGA	4620
TTGACCAAGC	CTGAGAGAGC	AACCTGCTGC	AATTGCTCAT	TTGTAGTAGG	CAGTTTTTTA	4680
GTAGTCTCTA	GGAGAGCAGC	AACTAAATCT	TCACTCAAAT	CATGTCGAGC	ATGATTGTAA	4740
AGATCTTTTA	TAAGGCTTTC	TAGGTTTGGT	TCTACCATCC	CTACCACCTC	CCTTATGGTT	4800
TAATAATGTT	TAATCAAAAT	AACCGTTGAA	CGATCCAATT	TCTTCACCAA	GCCTTGTAA	4860
AAAGCTTTCG	CTTCTAGGAA	GTCACTCAAT	GCATAGAGGG	TTTGGTGAGA	ATGGATATAA	4920
CGAGCGCAGA	CACCGATAGT	TGTTGATGGG	ACACCACCAT	TTTTCATGAT	AGCTGCACCT	4980
GCATCTGTTC	CGCCTTTTAC	ACAGTAGTAT	TGGTACTTGA	TACCAGCTTC	TTTACGCGTT	5040
GTCAAAAGGA	AATCTTTCAT	CCCTGGGAGA	AGCAAGTGAC	CTGGATCATG	GAACGAAATC	5100
AAGGTTCCAT	CTCCAATCTT	GCCTTGACCA	CGTAGACAT	CACCTGCTGG	TGAGCAATCA	5160
ACTGCGAGGA	AGACTTCTGG	GTCAAACTTG	GTTGTAGAGG	TATGAGCGCC	ACGCGAGCCA	5220
ACTTCTTCTT	GGAGCTTAGA	ACCCAGATAG	AGTTCAATTC	CGAGTTTTTG	ACCCGATAAA	5280

273

GCTTCAGCTA GCTGCTTAC CATGAGGACA CGTAGCGGT TATCCCAAGC TTTTGAGATG 5340  
 ATATTTTTTT CATTTGGCTGT CAAAAATGCA GAACATATCTG GTACAATGGT ATCACCAGGA 5400  
 CGGATGCCAA AACTTTCTGC CTCAGCCTTG TCCGAAAAC CACCATCAA AACGATATCG 5460  
 GCAATGGCTG GCATGGTTGG TCCCCCTTT CCACGAGTCA AATGCGGAGG AACGAAACCT 5520  
 GAATTCACAG GAATTTTCATG ACCATCACGA GTCAAGAGTT TGAAACGTG GCTGCTAACC 5580  
 ACCATGGGGT TCCAGCCACC GATTTCTACG ACACCGAAGG TACCATCTGG CTTGATTTCTG 5640  
 CTGACCATAA AACCAACTTC GTCCATATGA GAAGGACCA AGACGCGGG TGCATCCACA 5700  
 GCTTCTGAAT GTTTGATACC AAAAATACCA CCCAAGCCAT CTGTCAACC TCTATCCACA 5760  
 TCGCGTGTCA ACTTTTCACG AAGATAAGCA CGGACAGGCG CTTCATGACC TGAGACTGCA 5820  
 GCAAGTCTG TTACTTCTTT AATTTTIGAA AATTAATGTT TCAATTCAGT TCCTTCTTTT 5880  
 TTTTCATCAT TTTACCATT TTTATAGGAG AAGGATAGTG GGAAGGTGGA TTTCTAAGTT 5940  
 AGTATCTTAG TCTGCTCTA TCTTAGAAA GGATAGTATT CTCTTGCA TGATGCAAAA 6000  
 TCTAGTAAC ATTCGAAAT TAACTCGAAT ATTTATTTCC AAACAAAAA ACAATACACC 6060  
 ATCAAGATTG TTTGGATTT TCAATGAAAT TACAGAAAT AGTTGACTTC CCTTCTCTCT 6120  
 GTCTTTAAAT ATATAGTTGG TTGAGTTTGG AATAGTACGC TGTAGCTGCT AAAACATTTT 6180  
 TAGAAATTAA TTTGACTTTC CTAATAGAGT TGTTCATATC TTAATTCAT TTAATATAGT 6240  
 ACAGAACTAG AAAAGGAAA AATCATGACC AGG 6273

(2) INFORMATION FOR SEQ ID NO: 22:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 28171 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 22:

ACAACCTTTT TCAAAAACCT ACCTGGTAC GGAGATGTTT TGCTTCTGCT TATTATTTTC 60  
 GGTATATTC ATATCAATTT TGCTTTAACT CCTCTTGCTT TTTTCATTTTA TGCTAGTGGA 120  
 GGTCTTATTT TAGCTCTATT GTATCGCATG ACTAAAAATC TCTACTATCC AATACTAGTT 180  
 CATATCTCA TTAATATCAC TGCCTTCTGG GATGTGTGGT TGCTCCTATT TTGAGGAAGT 240  
 TAGCTTACTA AATAATGTC GGAACCTTCC GGCATTTTCT TTTTTCACAA ATAGTCAAG 300  
 TTTTCTCTTT CGATATTTGA GTGCTGTGTA TCCAGTTATT TTTTGAATTT GATTTTGAAA 360

274

ATAAGGTTGA CPTGAGAAAG GCAGATAGTG AAGATAGTTA AGAAGANTAG GATGTTCTTT	420
TTTCTCTTTT GGAAGAACTTC TAAATATGCG TATAATGAAA AGATAAAGAA GTTGGGGGTA	480
GAAGATGAAC ATTCAACAAT TACGCTATGT TGTGGCTATT GCCAATAGTG GTACTTTPTCG	540
TGAAGCTGCT GAAAGATGT ATGTTAGTCA GCCGAGCTCG TCTATTTCTG TTGCTGATT	600
GGAAAAGAG TTGGGCTTTA AGATTTTCCG TCGGACAGC TCAGGGACTT TCTTGACCG	660
TCGTGGGATG GAATTTTATG AAAAATCGCA AGAATTGGT AAAGGATTG ATATTTTCA	720
AAATCAGTAT GCCAATCTCG AAGAAGAAA AGATGAATTT TCTGTGCTA GCCAGCACTA	780
TGACTTCTTG CCACCAACTA TTACGGCCTT TTCAGAGCGC TATCCTGACT ATAAGAACTT	840
CCGTATTTTT GAATCAACTA CTGTCTAAAT ATTAGATGAA GTGGCGCAAG GGCATAGTGA	900
GATTGGGATT ATCTACCTCA ACAATCAAAA TAAAAAGGG ATTATGCAAC GGGTTGAAAA	960
ATTAGGCTCG GAGGTCATCG AATTGATTCC TTTCATACC CATATTTATC TCGGTGAGGG	1020
TCATCTTTTA GCCCAGAAAG AGGAATTAGT CATGAGGAT TTAGCGGATT TACCAACGGT	1080
TCGTTTCACT CAAGAGAAAG ACGAGTACCT TTATTTTCA GAGAACTTTG TCGATACCAAG	1140
CGTAGCTCA CAGATGTTTA ATGTGACAGA CCGTCCACC TTGAATGGTA TTTTGGAGCG	1200
GACGAGCGCT TATGCGACAG GTTCTGGATT TTTAGATAGT GACAGTGTTA ATGGCATTAC	1260
AGTTATTCTG CTCAAGGATA ACCTAGATAA CCGCATGGTC TATGTTAAAC GTGAAGAAGT	1320
GGAGCTTAGT CAAGCTGGGA CTCTCTCGT AGAAGTCATG CAAGAATATT TTGATCAAAA	1380
GAGGAATCA TGAAAAAAG AGCAATAGTG GCAGTCATTG TACTGCTTTT GATTGGGCTG	1440
GATCAGTTGG TCAATCTCTA TATCGTCCAG CAGATTTCCAC TGGGTGAAGT GCGCTCCTGG	1500
ATCCCCAATT TCGTTAGCTT GACCTACCTG CAAAATCGAG GTGCAGCCTT TTCTATCTTA	1560
CAAGATCAGC AGCTGTTATT CGCTGTCAAT ACTCTGGTTG TCGTGATAGG TGCCATTGG	1620
TATTTACATA AACACATGGA GGAATCATTC TGGATGGTCT TGGGTTTGAC TCTAATATTC	1680
GCGGTTGGTC TTGGAACACT TATTGACAGG GTCACTCAGG GCTTTGTGT GGATATGTTT	1740
CACCTTGACT TTATCAACTT TGCAATTTTC AATGTGCGAG ATAGCTATCT GACGGTTGGA	1800
GTGATTTTTT TATTGATTGC AATGCTAAAA GAGGAATAA ATGGAAATTA AATTTGAAC	1860
TGCTGGTCTG CGTTTGATA AGGCTTTGTC AGATTGTGCA GAATTATCAC GTAGTCTGCG	1920
GAATGAACAA ATTAAATCAG GCCAGTCTT GGTCAATGOT CAAGTCAAGA AAGCTAAATA	1980
CACAGTCCAA GAGGTTGATG TCGTCACTTA CCATGTGCCA GAACCAGAGG TATTAGAGTA	2040
TGTGGCTGAG GATCTTCCGC TAGAAATAGT CTACCAAGAT GAGGATGTGG CTGCTGTTAA	2100
CAAACTCTAG GGAATGGTTG TGCACCGAG TGCTGGTCAT ACCAGTGAA CCCTAGTAAA	2160



TGCCCTCATG	TATCATATTA	AGGACTTGTC	GGGTATCAAT	GGGGTTCTGC	GTCCAGGGAT	2220
TGTTACCGGT	ATTGATAAGG	ATACGTACAG	TCTTCTCATG	ATTGCTAAAA	ACGATGATGC	2280
GCATCTAGCA	CTTGCCCAAG	AACTCAAGGA	TAAAAAGTCT	CTCCGCAAA	ATTGGGCGAT	2340
TGTTCAAGGA	AATCTACCTA	ATGATCGTGG	TGTAATTGAA	GCGCCGATTG	GCCGAGTGGA	2400
AAAAGACCGT	AAGAAACAGG	CTGTAACCTG	TAAAGGGAAG	CCCTGCAGTA	CGCGTTTTC	2460
CGTCTTGGA	CGCTTTGGCG	ATTATAGCTT	AGTAGAGTTG	CAACTGGAGA	CAGGGCGCAC	2520
TCATCAAAATC	CGTGCCACA	TGGCTTATAT	CGGCATCCA	GTGCTGGTG	ATGAGTCTTA	2580
TGGTCTTCGC	AAGACTTTGA	AAGGACATGG	ACAATTTCTT	CATGCCAAGA	CTTTAGTCTT	2640
TACTCATCCG	AGAACAGGTA	AGACCTTGGA	ATTTAAAGCA	GATATCCCG	AGATTTTATA	2700
GGAAACCTTG	GAGAGATTGA	GAAGTAAGA	ATGAAAAGA	AATTAAGT	TTTAGCACTT	2760
GTAGGCGCTT	TTTTAGGTTT	GTCAAGGTAT	GGGAATGTTT	AGGCTCAAGA	AAGTTCAGGA	2820
AATAAAATCC	ACTTTATCAA	TGTTCAAGAA	GGTGCGAGTG	ATGCCGATT	TCTTGAAAGC	2880
AATGGACATT	TGCGCATGGT	GGATACAGGA	GAAGATTATG	ATTTCCGAGA	TGGAAGTGAT	2940
TCTCGCTATC	CATGGAGAGA	AGGAATTGAA	ACGTCCTATA	AGCATGTCTT	AACAGACCGT	3000
GTCTTCTGTC	GTTTGAAAGA	ATTGGGTGTC	CAAAAACCTG	ATTTTATTTT	GGTGACCAT	3060
ACCCACAGTG	ATCATATTGG	AAATGTTGAT	GAATTACTGT	CTACCTATCC	AGTTGACCCA	3120
GTCTATCTTA	AGAAATATAG	TGATAGTCGT	ATTACTAATT	CTGAACGCTT	ATGGGATAAT	3180
CTGTATGGCT	ATGATAAGGT	TTTACAGACT	GCTGCAGAAA	AAGGTGTTTC	AGTTATTCAA	3240
AATATCACAC	AAGGGGATGC	TCATTTTCAG	TTTGGGACA	TGGATATTCA	GCTCTATAAT	3300
TATGAAAAATG	AAACTGATTG	ATCGGGTGAA	TTAAAGAAAA	TTTGGGATGA	CAATTCCAAT	3360
TCCTTGATTA	CGGTGGTGAA	AGTCAATGGC	AAGAAAAATT	ACCTTGGGGG	CGATTTAGAT	3420
AATGTTCAATG	GAGCAGAGGA	CAAGTATGGT	CCTCTCATTG	GAAAAGTTGA	TTTGATGAAG	3480
TTTAATCATC	ACCATGATAC	CAACAAATCA	AATACCAAGG	ATTTTATATA	AAATTTGAGT	3540
CCGAGTTTGA	TTGTTCAAAC	TTCCGATAGT	CTACCTTGGA	AAATGGTGT	TGATAGTGAG	3600
TATGTTAATT	GGCTCAAGGA	ACGAGGAATT	GAGGAATCA	ACGACGCCAG	CAAAAGCTAT	3660
GATGCAACAG	TTTTTGATAT	TCGAAAAGAC	GGTTTGTGTA	ATATTTTCAAC	ATCCTACAAG	3720
CCGATTCCAA	GTTTTCAAGC	TGTTGGCAT	AAGAGTGCA	ATGGGAACGT	GTGGTATCAA	3780
GCGCTTGATT	CTACAGGAGA	GTATGCTGTG	GGTTGGAATG	AAATCGAAGG	TGAATGGTAT	3840
TACTTTAAAC	AAACGGGTAT	CTTGTTACAG	AATCAATGGA	AAAAATGGA	CAATCATTTG	3900

		276	
TTCTATTGGA	CAGACTCTGG	TGCTCTGCT	AAAAATGGA AGAAAATCGA TGGAACTCGG 3960
TATTATTTTA	ACAAAGAAAA	CCAGATGGAA	ATGGGTGGA TTCAAGATAA AGAGCAGTGG 4020
TATTATTTGG	ATGTTGATGG	TTCTATGAAG	ACAGGATGGC TTCAATATAT GGGGCAATGG 4080
TATTACTTTG	CTCCATCAGG	GGAAATGAAA	ATGGGCTGGG TAAAAGATAA AGAAACCTGG 4140
TACTATATGG	ATTCTACTGG	TGTCATGAAG	ACAGGTGAGA TAGAAGTTGC TGGTCAACAT 4200
TATTATCTGG	AAGATTCAGG	AGCTATGAAG	CAAGGCTGGC ATAAAAAGGC AAATGATTTGG 4260
TATTTCTACA	AGACAGACGG	TTACAGAGCT	GTGGGTTGGA TCAAGGACAA GGATAAATGG 4320
TACTTCCTGA	AAGAAAATGG	TCAATTACTT	GTGAACGOTA AGACACCAGA AGGTTATACT 4380
GTGGATTCAA	GTGGTGCCCTG	GTTAGTGGAT	GTTTCGATCG AGAAATCTGC TACAACTAAA 4440
ACTACAAATC	ATTCAGAAAT	AAAAGAATCC	AAAGAAATAG TGA AAAAGGA TCTTGAAAT 4500
AAAGAAACGA	GTCAACATGA	AAGTGTACAA	AAATTTTCAA CTAGTCAAGA TTTGACATCC 4560
TCAACTTCAC	AAAGCTCTGA	AACGAGTGTA	AACAAATCGG AATCAGAACTA GTAGTAGAAA 4620
AGAAAGTTT	AGGGCCCTCT	TTTTCTCTATC	AATCTTTTTC TATTTCTCTG TATTCATGTT 4680
ATAATGGATA	AATATGAATA	ATCGGAGTGA	GACTATGAAA TCAAAACGGA TTCTCTTTAA 4740
GGTGGGTACT	TCTTCTCTGA	CAAAATGAGGA	TGGAACTTTA TCACGTAGTA AGGTAAAGGA 4800
TATTACCCAG	CAGTTGGCTA	TGCTGCACGA	GGCTGGTCTAT GAGTTGATTT TGGTGTCTTC 4860
AGGTGCCATT	GCGGCTGGTT	TTGGAGCCTT	AGGATTTTAA AAGCGTCCGA CTAAGATTGC 4920
TGATAAACAG	GCTTCAGCAG	CGGTAGGGCA	GGGGCTTTTG TTGGAAGAAAT ATACAACCAA 4980
TCTTCTCTTG	CGTCAAAATCG	TTTCTGCACA	AATCTTCTGT ACCCAAGATG ACTTTGTGGA 5040
TAAAGCTCGT	TATAAAAAATG	CCCATCAGGC	TTTGTCTGTT TTGCTCAACC GTGGGGCAAT 5100
TCCTATCATC	AATGAGAAATG	ATAGTGTGCT	TATTGATGAG CTC AAGGTTG GGGACAATGA 5160
CACCTCTAAGT	GCTCAAGTAG	CGGCGATGGT	CCAAGCAGAC CTTTATGTTT TCTTTGACAGA 5220
TGTGGACGGT	CTCTATACTG	GAAATCCTAA	TTTCAGATCCA AGAGCCAAC GCTTTGAGAG 5280
AATCGAGACC	ATCAATCTGT	AGATTATTGA	TATGGCTGGT GGAGCTGGTT COTCAACCGG 5340
AACCTGGGGT	ATGTTAAACCA	AAATCAAGGC	TGCAACTATC GCGACGGAAT CAGGAGTTCC 5400
TGTTTATATC	TGCTCATCCT	TGAAATCAGA	TTCCATGATT GAGGCGGAG AGGAGACCGA 5460
GATGCTCTCT	TACTTTGTTG	CTCAAGAGAA	GGGGCTTCTG ACCCAGAAAC AATGCTTTGC 5520
CTTCTATGCT	CAGAGTCAAG	GTTCTATTTG	GTTTGATAAA GGGGCTGGG AAGCTCTCTC 5580
TCAATATGGA	AAGAGCTTTC	TCTTTATCTGG	TATCGTTGAA GCAGAAGGAG TCTTTTCTTA 5640
CGGTGATATC	GTGACGATAT	TTGACAAGGA	AAGTGGAAAA TCACTTGAAA AAGGACCGGT 5700

GCAATTGGGA GCATCTGCTT TGGAGGATAT GTGCGTTCT CAAAAGCCA AGGGTGTCTT 5760  
 GATTTACCGT GACGACTGGA TTTCATPAC TCCTGAATC CAACTACTTT TTACAGAATT 5820  
 TTAGAGGTAA ACTATGGTGA GTAGACAAGA ACAATTGAA CAGGTACAGG CTGTTAAAAA 5880  
 ATCGATTAA ACAGCTAGTG AAGAAGTCAA AAACCAAGCC TTGCTAGCCA TGGCTGATCA 5940  
 CTTAGTGGCT GCTACTGAGG AAATTTTAGC GGCTAATGCC CTCGATATGG CAGCGGCTAA 6000  
 GGGGAAAATC TCAGATGTGA TGTGGATCG TCTTATTTG GATGCAGATC GTATAGAAGC 6060  
 GATGGCAGA GGAATTCGTG AAGTGGTTG CTTACCAGAT CCAATCGGTG AAGTTTATGA 6120  
 AACAACTCAG CTTGAAAATG GTTTGGTTAT CACAAAAAA CUGTAGCTA TGGGTGTCTT 6180  
 CGGTATTATC TATGAAAGCC GTCCAAATGT GACGTCTGAT GCGGCTGCTT TGAATCTTAA 6240  
 GAGTGGAAAT GCGGTGTCTT TCTGATGTG TAAGGATGCC TATCAAAACA CCCATGCCAT 6300  
 TGTCTACAGC TTGAAGAAGG GCTTGGAGAC GACTACTATT CATCCAAATG TGATTCAACT 6360  
 GGTGGAGGAT ACTAGCCGTG AAAGTAGTTA TGCTATGATG AAGGCCAAGG GCTATCTAGA 6420  
 CCTCTCATTT CCTCTGGAG GAGTGGCTT GATCAATGCA GTGGTTGAGA ATGCGATTGT 6480  
 ACCTGTTATC GAGACAGGGA CTGGGATTGT CCATGTCTAT GTGGATAAGG ATGCAGACGA 6540  
 AGCAAGGCG CTGTCTATCA TCAACAATGC TAAACCAGT CTTCTTCTG TTTGTARTGC 6600  
 CATGGAGGTT CTGCTGTTT ATGAAAAA GGCAGCAAGC TTCCTTCTC GCTTGGAGCA 6660  
 AGTGTGGTT GCAGAGCGTA AGGAAGCTGG ACTGGAACCA ATTCAAATCC GCCTAGATAG 6720  
 CAAGCAAGC CAGTTGTTT CAGTCAAGC AGCTGAGACC CAGACCTTG ACACCGAGTT 6780  
 TTTAGACTAT GTCTTGTCTG TTAAGGTTGT GAGCAGTTTA GAAGAAGCGG TTGCGCACAT 6840  
 TGAATCCAC AGCACCATC ATTCGGATGC TATTGTGAGG GAAAAAGCTG AAGCTGCAGC 6900  
 ATACTTTACA GATCAAGTGG ACTCTGCAGC GGTGTATGTT AATGCCCAA CTCGTTTAC 6960  
 AGATGGAGGA CAATTGGTC TTGGTGTGA AATGGGATT TCTACTCAGA AATTCACCC 7020  
 GCGTGTGCC ATGGGCTTGA AAGAGTTGAC CAGTACAAG TATGTGTTG CCGGTGATGG 7080  
 GCAGATAAGG GAGTAAGAGA TGAAGATTG ACTTATCGGT TTGGGGAATA TGGTGCTAG 7140  
 CTTGGCAAAA TCTGTCTTGC AGACTAGGAC GTCAGATGAG ATTCTCCTTG CCAATCGTAG 7200  
 TCAAGCTAAG GTAGATGCTT TCATTGCGAGA CTTTGGTGT CAGGCTTCCA GCAATGAAGA 7260  
 AATGTTTCCA GAAGCAGATG TGATTTTCT AGGAGTTAAG CCTGCTCAGT TTTCTGAAT 7320  
 GCTTCTCAA TACCAGACCA TCCTTGAAAA AAGAGAAAGT CTCCTTTTGA TTTCGATGGC 7380  
 AGCTGGATTG ACCTTAGAAA AACTAGCAAG TCTTATCCCA AGTCAACACC GATTATTTCG 7440

278

TATGATGCCT AATACCCCTG CTTCATCCG GCAAGGAGTG ATTAGTTATG CCTTGCTCC	7500
TAAATGCAGG GCTGAGGACA GTGAGCTCTT TATCAGCTT TTAGCCAGG CTGCTCTCTT	7560
GTTTGAAC TAAGAAAGTT TAATCGATCG AGCGACAGGT CTTCAGAGTT GTGGACCAGC	7620
CTTTGCTTAT CTTTTTATCG AGGCTTTGCC AGATGCAGGT GTTCAGACAG GATTACCAAG	7680
AGAAATAGCA TTGAAATGG CAGCACAAC TGTGGTAGGA GCTGGGCAAT TGGTCCCTGA	7740
AAGTCAGCAA CATCTGGAG TATTGAAAGA CCAAGTCTGT AGCCCAAGCG GTTCGACTAT	7800
CGCTGGTGA GCAAGCCTAG AAGCGCATCG TTTCCGAGGA ACAGTCATGG ATGCAGTTCA	7860
TCAGGCTAC AAACGAACAC AAGAACTAGG TAAATAAGAG GTAGTTTGA CTGCTCTCTT	7920
TATGTTGGT GAAATGAGAA GACACAAAAA GATTGTGACA AACCCCTATT TTTTGTATAG	7980
AATAGAAATA GTAAAAAGA AATGAGTTAG ACATGTCAA AGGATTTTGA GTCTCTCTTG	8040
AGGGACCAGA GGGAGCAGCG AAGACCACTG TTTTAGAGGC TCTGCTACCA ATTTTAGAGG	8100
AAAAGGAGT AGAGGTGTTG ACGACCCGTG AACCTGGCGG AGTCTGTATT GGGGAGAGA	8160
TTCCGGAGAT GATTTTGGAT CCAAGTCATA CTCAGATGGA TGCTAAAAACA GAGCTACTTC	8220
TCTATATTGC CAGTCGCAGA CAGCATTTGG TGGAAAAAGT TCTTCCAGCC CTGGAAGCTG	8280
GCAAGTTGGT CATCATGGAT GTTTTTATCG ATAGTTCTGT TGCCTATCAG GGATTGGTC	8340
GTGGCTTAGA TATTGAAGCC ATTGACTGGC TCAATCAGTT TGCAGACAGT GGCTCAAAC	8400
CCGATTTGAC ACTCTATTTT GACATCGAGG TGAAGAAGG GCTGGCTGTG ATTGCTGCTA	8460
ATAGTGACCG CGAGGTTAAT GTTTTGGATT TGGAAAGGTT GGAATTCAGT AAAAAAGTTC	8520
GTCAAGGCTA CCTTCTCTTT CTGGATAAAG AGGGAATCG CATTTGCAAG ATGATGCTTA	8580
GTCTCCCTTT GGAGCAAGTT GTGGAAACTA CCAAGGCTGT CTGTGTTGAC GGAATGGCT	8640
TGGCCAAATG AAACAAGATC AACTAAAGGC TTGGCAACCA GCTCAGTTTG ACCGTTTTGT	8700
CCGATCTTTA GAACAAGACC AGCTCAATCA CGCTATCTC TTTTCAGGTT TCTTTGAAAG	8760
CTTGAAATG GCGCAATTTT TAGCTAAGAG CCTCTTTTGT ACGGATAAAG TTGGCGTCTT	8820
ACCATGTGAG AAATGCCGAA GTTGCAAGCT GATTGAACAG GGAGAAATTC CCGATGTAC	8880
CTTGATTAAG CCAAGTTAAT AGGTCATTAA GACGGAAAGC ATTGAGAAAT TGGTGGGTCA	8940
GTTTTCTCAA GCAGGGATTG AAGGCCAGCA ACAGGTCTTT ATCATCGAGC AAGCGGATAA	9000
AATGCATCCC AAGCAGCCA ATTCTCTGCT CAAGGTCACT GAAGAAGCC AGAGTGAGAT	9060
TTATATTTTC TTCTTGACTA GCGATGAGGA AAAGATGTTA CCGACAAATC GAAGTCGGAC	9120
TCAGATCTTC CACTTTAAAA AGCAAGAAGA AAACCTTATC TTAATCTTAG AACAAATGGG	9180
ACTTGTAAAG AAAAAAGCA CTCCTTTAGC TAAGTTTAGT CAATGCGAG CTGAAGCAGA	9240

AAAGTTGGCT	AATCAGGCAA	GTTTTGGAC	CTTGGTCGAT	GAAAGTAAC	GCCGCTGAC	9300
TTGGTTAGTA	GCTAAGAAA	AAGAAAGTTA	TCTACAGGTT	GCCAAATTAG	CCAACTTGGC	9360
AGATGATAAG	GAAAAACAGG	ATCAGGTTTT	ACGATTCTCT	GAACTCTCT	GTGGGCAGGA	9420
CCTCTTGCAG	GTAAGAGTAA	GAGTGATTCT	ACAAGATTTA	CTAGAAAGCTA	GAAAAATGTG	9480
GCAAGCTAAT	GTCAGCTTTC	AAATGCCAT	GGAATATCTG	GTCTTGAAAG	AAATATAAAC	9540
TCAAAAATGA	ATGATAAAGA	AAGGAAAGGG	CTGTTTTATG	GACAAAAAAG	AATTATTTGA	9600
CQCGCTGGAT	GATTTTCCC	AACAAATTAT	GGAACCTTA	GCCGATGTGG	AAGGCATCAA	9660
GAAAAATCTC	AAGAGCTTGG	TAGAGGAAA	TACAGCTCTT	CGCTTGGAAA	ATAGTAAGTT	9720
CCGAGAACGC	TTGGGTGAGG	TGGAAGCAGA	TGCTCTGTGC	AAGGCCAAGC	ATGTTCGTGA	9780
AAGTGTCCGT	CGCAATTACC	GTGATGGATT	TACAGTATGT	AATGATTTTT	ATGGACAAGG	9840
TCGAGAGCAG	GACGAGGAAT	GTATGTTTTG	TGACGAGTTG	CTATACAGGG	AGTAGGCATG	9900
CAGATTCAAA	AAAGTTTAA	GGGCGAGTCT	CCCTATGGCA	AGCTGTATCT	AGTGCCAAAG	9960
CCGATTGGCA	ATCTAGATGA	TATGACTTTT	CGTGCTATCC	AGACCTTGAA	AGAAGTGGAC	10020
TGGATIGCTG	CTGAGGATAC	GCGCAANTACA	GGGCTTTTGC	TCAAGCATTT	TGACATTTCC	10080
ACCAAGCAGA	TCAGTTTTCA	TGAGCACAAAT	GCCAAGGAAA	AAATTCCTGA	TTTGTATTGGT	10140
TTCTTTGAAJG	CAGGGCAAG	TATTGCTCAG	GTCTCTGATG	CCGGTTTGCC	TAGCATTTCA	10200
GACCCITGGT	ATGATTTAGT	TAAGGCAGCT	ATTGAGGAAG	AAATTCAGT	TGTGACAGTT	10260
CCAGGTGCCCT	CTGCAGGAAT	TTCTGCCTTG	ATTGCCAGTG	GTTTAGCGCC	ACAGCCACAT	10320
ATCTTTTACG	GTTTTTFACC	GAGAAAAATCA	GGTCAGCAGA	AGCAATTTTT	TGGCTTGAAA	10380
AAAGATTATC	CTGAAACACA	GATTTTTTAT	GAATCACCTC	ATCGTGTAGC	AGACACGTTG	10440
GAAAAATATGT	TAGAAAGCTA	CGGTGACCGC	TCCGTGTCT	TGGTCAGGGA	ATTGACCAAA	10500
ATCTATGAAG	AATACCAACG	AGGTACTATC	TCTGAGTTAT	TAGAAAGCAT	TGCTGAAAGG	10560
CCACTCAAGG	CGCAATGTCT	TCTCATTTGT	GAGGGTGCCA	GTCAAGGTGT	GGAGGAAAAG	10620
GACGAGGAAG	ACTTGTGTGT	AGAAATTCAA	ACCCGCATCC	AGCAAGGTGT	GAAGAAAAAC	10680
CAAGCTATCA	AGGAAGTCGC	TAAGATTTAC	CAGTGGAAAT	AAAGTCAGCT	CTACGCTGCC	10740
TACCAAGCAT	GGGAAGAAA	ACAAATAAAG	GAGCAGGAT	GTAATAAATC	TGTCTGTTC	10800
TGTTTAACCTT	AATTAGTGAT	GATAATATAA	AGATGTATCA	CTTGATATAG	AAGCTTTGGT	10860
ATTAAAGTTT	TTATTAAGCC	CATACGGAAT	ACCGATGGTT	GGAGCAGCAG	TTATAGCGTT	10920
CTTAGAAGGT	ATAAATAGAA	AAATAAGGTC	ATTTTAATTC	AAAGGATTGA	TAAATCAGAA	10980

289

AGAAGGTGAT	TTTTGCGAA	CATACGAAA	TAAAGSAGAA	CTAAAAGCTG	AGATAGAGAA	11040
AACATTTGAG	AAATATATTT	TAGAAATTTGA	TAATATTCCA	GAATAATTAA	AAGATAAGAG	11100
AGCTGATGAA	GTTGACAGAA	CTCCAGCAGA	AAACCTTGCT	TATCAGGTTG	GTTGGACCAA	11160
CTTGCTCTTT	AAATGGGAAG	AAGATGAAAG	AAAGGGGCTT	CAAGTAAAAA	CACCATCGGA	11220
TAAATTTAA	TGGAATCAAC	TGGTGAATT	ATATCAGTGG	TTACAGATA	CTTACGCTCA	11280
TTTATCTCTG	CAAGAGTTGA	AAGCAAAATT	AAATGAAAT	ATTAAATTCTA	TCTCTGCAAT	11340
GATTGATTCG	TTGAGTGAGG	AAGAATTATT	TGAACCGCAT	ATGAGAAAGT	GGCTGATGA	11400
AGCGACTAAA	ACAGCGACTT	GGGAAGTGT	TAAGTTTATT	CATGTAAATA	CGGTTGCACC	11460
TTTTCGAACT	TTCAGAACTA	AAATCAGAAA	ATGGAAGAAG	ATAGTATTAT	AAATTATATT	11520
TTTAACTTTA	AAAAATTTCA	TAAAAATGGT	TACCAAAGGC	GATAGAAGAA	AAACTATCGT	11580
CTTTTCTTT	GCAAAATTTT	AAGAAAGGAG	GTGATCTTGC	ATGGACTTTG	AATATTTTTA	11640
TAACAGAGAA	CGGAAAGAT	TTAATTCTTT	AAAAGTACC	GAGATATTAG	TTGATAGAGA	11700
AGAAATTCGG	GGCTTATCAG	CAGAAGCAAT	TATCCTTTAT	TCCATACTTC	TTAAACAGAC	11760
AGGAATGTCA	TTTAAAGAATA	ACTGGATAGA	CAAGGAAGGC	AGAGTATTTA	TCTATTTTAC	11820
TGTCGAAGAA	ATTATGAAAA	GAAGAATAT	CTCAAAGCCA	ACTGCCNTAA	AAACATTAGA	11880
TGAGCTTGAT	GTAAGAAAAGG	AATAGGACTG	ATCGAAAGAG	TAAGGCTTGG	ACTTGGTAA	11940
CCGAACATCA	TTTATGTTAA	AGACTTTIATG	AGTATATTTT	AGGTAAAAA	AAATGACTTA	12000
CAGAAGTCAA	AAAACTTAAC	TTCAGAAGTA	AAAGATTTTA	ACCTCAGAG	TAAAGAAAA	12060
GAACITCAAG	AGGTTAAGAA	CCTTGACTCT	AACTATATAG	AGAAATAATA	GAGTAAGTAT	12120
AGTAAGAGAG	AATATAGTTT	TGGTGAAAAC	GGACTTGGA	CATTTCAAAA	TGTGTTTTTA	12180
GCTGCTGAAG	ATATATCGGA	TTTACAAATC	ATAATGAAC	CACAGCTTGA	GAATTACATT	12240
AGACTTCTTG	CAAACTAGA	ATCCTAGTTC	ATGATTGATA	ATGCCAGCA	TCAAATTCAT	12300
TCGTATTCG	AAGCGTTTAC	GATGATTTCG	ATAGATTGTT	GAAACATTT	TAAACGTTTT	12360
TACTTTGGCA	AAGATGTTCT	CAATCTTGCT	TCTCTCCTTG	GATAGCGCAT	GGTTACAGGC	12420
TTTATCTTCA	GCTGTTAGCG	GCTTGAGTTT	GCTGGAATTA	CGTGAGTTT	GTACTTGGAG	12480
ATTATCTTTC	ATGAGCCCTT	GATAACCACT	GTGAGACAAG	ATTTTACCAG	CTGTCCGAT	12540
ATTTCTGCGA	CTCAATTTGA	ACAACCTTCA	ATCACGACAA	TAGTTCACAG	CGATATCCAA	12600
AGAAACAATT	CTCCCTTGAC	TTGTGACAAT	CGCTTGAGCC	TTTATACGCT	GAAATTTCTT	12660
TTTACCAGAA	TGATTCGCTA	ATTCTTTTTT	TAGGGCGATT	GATTTTACT	TCCGTGCGAT	12720
CAATCATTAC	CGTGTCTCA	GAACTGAGAG	GAGTTCTTGA	AATCGTAACA	CCACTTTGAA	12780

CAAGAGTTAC TTCAACCCAT TGGCTCCGAC GGATTAAGTT GCTTTCOTGA ATACCAAAAT 12840  
 CAGCCGCAAT TTGTTCAATA GTTCGATATT CTCGCACATA TTGAAGAGTG GCCATAAGAA 12900  
 GGTCTTCTAG GCTTAATTTA GGTTCGCTC CACCTTTTGC GTGTTTAAGT TGATAAGCTG 12960  
 TTTTAAATAC AGCTAATATC TCTTCAAAAG TCGTGGCCTG AACACCAACA AGACGCTTAA 13020  
 ATCGTGATC AGTTAGTTGT TTACTTGGCT CATCATTCAT AGACTACTA TACCATATTT 13080  
 TGTITCCGAG GAAGTCTATT GGAAAGTAAG AAATATTGAA GCTGAGGCTA TTAGAAGAAA 13140  
 TTGTAGAGGT GGTGCTATTT TTTCAGGTAA AATAAAATAT CACGAAGATT CACAGTTTAA 13200  
 AGGAGATCAC TATGTTGAAT GTTATGCTGT TTTAGATAAT ACGTTATAG CAAGAGATAG 13260  
 AATAACAGTC CCTATCGATC CGTTATGTGG AAAAGATTTT ATAGAGTAGC ATATAATTGA 13320  
 TTCTTAAC TGACTACTAC TATCTCTTTA CATCAAGAAA ATGACTAAAC AGGGAAGTTT 13380  
 GCCTTCTTCC CTTTTTTTGT TATACTAGTA GAAGAAAAAA TTAGAAGAT TGTGGGTGT 13440  
 CAACAGCCCC AGTGGGGTGT TTTAATATGG ACTTAGGTCC CACCCAAAGA GGTATTAGTG 13500  
 TCGTCTCTCA ATCTATATAT AATGTTATCG GTGCTGGTTT GGCAGGTCTT GAAGCAGCTT 13560  
 ACCAAATCGC AGAGCGTGTG ATTCAGTTAA AACTATATGA AATCGGTGGT GTCAAGTCTA 13620  
 CACCCAGCA TAAACAGAC AATTTTGCTG AGTTGGTTTG TTCCAATTCT TTGCGTGGG 13680  
 ATGCTTTGAC AAATGCAAGT GGTCTTCTCA AGGAAGAAAT GCGTGGCTTG GGTCTGTTA 13740  
 TCTTGGAA TCCTGAGGCT ACACGTGTTT CTGCAAGTGG TGCCCTTGCA GTGACCGTG 13800  
 ATGCTTTCTC TCAANTGGTG ACCGAAAAAG TTGCCAACCA CCCCTTGATT GAAGTGGTTC 13860  
 GTGATGAAT TACAGAAITG CCGACAGATG TTATTACGGT TATCGCTACT GGTCTTTTGA 13920  
 CAAGTAGTGC CTGCGCTGAA AAGATTATG CTCTTAATGA CGGTGCTGGT TTTTATTTCT 13980  
 ACGATGCGGC AGCGCTTATT ATCGATGTCA ACACATATGA TATGAGCAAG GTCTACCTCA 14040  
 AATCACGTTA TGATAAGGA GAAGCGGCTT ACCTCAATGC CCCTATGACC AAGCAAGAT 14100  
 TTATGATTT CCATGAAGCT TTGCTCAATG CAGAAGAAGC ACCCTTATG TCTTTTGAAA 14160  
 AAGAAAAGTA CTTTGAAGGA TGTATGCTTA TCGAAGTCAT GGCCAAACGT GGCATTAAAA 14220  
 CTATGCTTTA TGGCCCTATG AAGCCAGTCG GTCTTGAGTA CCCAGACGAC TATACAGGAC 14280  
 CTCGTGATGG AGAATTTAAA ACACCTTATG CGGTGTGCA ACTTCGTGAG GATTAATGAC 14340  
 CTGTAGCCCT CTACAATATT GTTGGTTTCC AGACCCACCT CAATGGGGA GAACAAAAGC 14400  
 GTGCTTCCA AATGATTCGG GGTCTTGAAA ATGCGAGATT TGTCCGTTAT GGTGTGATGC 14460  
 ATCGCAATTC TTACATGGAT TCACCAAATC TTCTTGAGCA GACTTACCGT TCTAAGAAAC 14520

282

AACCAATCT CTCTTTGCT GGTCAAATGA CGGCTGTGA AGGCTATGTT GAGTCGGGG	14580
CTTCAGGCTT AGTTGCGGA ATTAACGCGAG CTCGTCTCTT CAAGGAAGAA AGCGAGGCTA	14640
TTTTCCCGA GACGACAGCG ATTGGAAAGCT TAGCTCATTA CATTACCCAT GCCGACAGCA	14700
AACATTTCCA ACCAATGAAT GTCAATTTTG GGATCATCAA GGAGTTGGA GCGAGCGTA	14760
TCCGTGATAA GAAGGCTCGT TATGAAAAAA TTGCAGAGCG TGCCCTTGCC GACTTAGAGG	14820
AATTTTGTAC TGTCTAATTT TTTTGAAGA ATTGCTCATG ATACTATAAA AATCTTAGAA	14880
ATTGTGATAA AATAGGTAGG ATGAAGAAG GAGAGTGAAA ATGGCGAATC CCAAGTATAA	14940
ACGTATTTTA ATCAAGTTAT CAGGTGAAGC CCTTGCCGCT GAACGTGGCG TAGGGATTGA	15000
TATCCAAACA GTTCAACAA TCGCAAGA GATTCAAGAA GTTCATAGCT TAGGTATCGA	15060
AATTGCCCTT GTTATCGGTG GAGGAAATCT CTGGCCTGGA GAACCTGCAG CAGAAGCAGG	15120
TATGGACCGT GTTCAGGCAG ATTAACACAG AATGCTTGGG ACTGTTATGA ATGCTCTTGT	15180
GATGGCAGAT TCATGTCAAC AAGTTGCGGT TGATAAGCGT GTACAAACAG CTATTGCCAT	15240
GCAACAAGTG CGAGAGCCTT ATGTCCTGG ACCTGCCCTT CGTCACCTTG AAAAAGGCGG	15300
TATGCTTATC TTGTGTGCTG GAATTTGTTT ACCTTACTTC TCGACAGATA CAACAGCGGC	15360
CCTTGCTGCA GCTCAAAATC AAGCAGATGC CATCTCATG GCTAAAAATG GTGTCGATGG	15420
TGTTTCAAT GCGATCCTA AGAAGATTA GACAGCTGTT AAGTTTGAG AATTGACCCA	15480
CCGTGACGTT ATCAATTAAG GTCTTCGTAT CATGGACTCA ACAGCTCAA CCCTTCAAT	15540
GGACAACGAC ATTGACTTGG TTGTATTCAA CATGAACCAA CCAGGCAACA TCAACGTGT	15600
CGTATTTGGT GAAAATATCG GAACAACAGT TTCAAATAT ATCGAAGAAA AGGAATAGAA	15660
AAGAATATGG CTAACGCAAT TATTGAAAA GCTAAAAGAGA GAATGACCCA GTCTCACCA	15720
TCACATTGCTC TGAATTTGG TGGTATCCGT GCTGGTGGT CCAATGCAAG CTTCCTTGAC	15780
CGTGTACATG TAGAATACTA TGGAGTCGAA ACTCCTCTTA ACCAAATCGC TTCAATTACG	15840
ATTCCAGAAG CGCGTGTITT GTTGGTAACA CCAATTTGACA AGTCTTCATT GAAAGACATC	15900
GAACGTGGCT TGAACGCTTC TGATATTGGT ATCACACCGT CTAATGACGG TCTGTGAT	15960
CGCTTGGTTA TCCAGCTCT TACAGAAGAA ACTCGTCTG ACCTTGCTTA AGAAGTGAG	16020
AAGGTGCGCG AAAATGTCAA AGTGGCTGC CGCAATATCC CGACGGATG TATGGACGAA	16080
GCTAAGAAAC GAGAAAAAGC AAAAGAAATC ACTGAAGAC AATTGAAGAC TCTTGAAAAA	16140
GACATTCAAA AAGTAACAGA CGATGCTGTT AAACACATCG ACACATGAC TGCTAACAAA	16200
GAGAAAGAAC TTTTGGAGT CTAATAATTA ACAGAAAAAC TCAGTTGGCA TTGCTGGCTG	16260
AGTTTTATTC GAAAGAAAGA AATATGAATA CAATCTTGC AAGTTTATTC GTTGGACTGA	16320



TCATCGATGA	AAACGACCGT	TTTACTTTG	TGCAAAAGGA	TGGTCAAAAC	TATGCTCTTG	16380
CTAAGGAAGA	AGGCCAACAT	ACAGTAGGGG	ATACGGTCAA	AGGTTTGTGA	TACACGGATA	16440
TGAAGCAAAA	ACTCGGCTTG	ACAACCTTAG	AAGTGACTGC	CACCTCAGGAC	CAATTTGGTT	16500
GGGAGCTGT	CACAGAGGTT	CGTAAGGACT	TGGGTGTCTT	TGTGGATACA	GGCCTTCCTG	16560
ACAAGGAAAT	COTTGTGTCA	CTCGATATTC	TCCTTGAGCT	CAAGGAACCT	TGGCTTAGAA	16620
AGGGCGACCA	ACTCTACATC	CGTCTTGAAG	TGGATAAGAA	AGACCGATC	TGGGCGCTCT	16680
TGGCTTATCA	AGAAGACTTC	CAACGCTCTG	CTCGTCTCTG	CTACAACAAC	ATGCAGAAC	16740
AAACTGTGGC	AGCCATTGTT	TACCGTCTCA	AGCTGTCAAG	AACTTTTGTT	TACCTACCAG	16800
AAATTAATAT	GCTTGTTTTT	ATTATCTCTA	GCGAGCGTTA	CGCAGAGCCA	CGTTTGGGCT	16860
AGTATTAGA	TGCGGCGGTT	ATTGGTTTCC	GTGAAGTGGG	CGGCACTCTG	AACCTCTCCC	16920
TCAAACACAG	CTCCTTTGAA	ATGTTGGAAA	ACGATGCTCA	GATGATTTTG	ACTTATTTGG	16980
AAAGCAATGG	CGGTTTCATG	ACCTTAAATG	ACAAATCATC	TCCAGACGAC	ATCAAGGCCAA	17040
CCTTTGGCAT	TTCTAAAGGT	CAGTTCAAGA	AAGCTTTAGG	TGGTCTTATG	AAGGCTGGTA	17100
AAATCAAGCA	GGACCACTTT	GGGACAGAGT	TGATTTTAGG	AGGCTTATGA	GAJAATCATT	17160
TTACACTTGG	CTCATGACCG	AGCGCAATCC	TAAAAGTAAC	AGTCCCAAAG	CAATTTTGCG	17220
AGACCTCGCT	TTTGAAGAGT	CAGCCTTTCC	AAAACACACA	GATGATTTTG	ATGAGGTGAG	17280
TGCTTTTTTG	GAGGAGCATG	CCAATTTCTC	TTTTAACCTA	GGAGATTTTG	ACAGCATTTG	17340
GCAGGAATAT	CTAGAACAAT	AGCATTTATG	CATPGGGTTT	GGGCTAGTAA	TTTCTCCATC	17400
CCTCTGCTAT	AATAAAAAGA	AATAAAAAGG	TPAGAGAGGT	TCTTTATTTG	AAGGAACATT	17460
CAATAGACAT	TCAACTGAGT	CATCCAGATG	ACCTGTTTCA	TCTTTTGGT	TCCAATGAAC	17520
GCCATCTTCG	TTTGATGGAA	GAAGAGCTTG	ATGTTGTGAT	TCATGCTCGT	ACGGAGATTG	17580
TCCAGGTTTT	GGGAGAAGAG	TCTGCTGTGT	AGGAAGCCCG	TCAAGTTATT	CAGGCTTTGA	17640
TGGTCTTGTT	AAATCOTGGG	ATGACCGTTG	GTACGCCAGA	TGTAGTCACT	GGGATTAGCA	17700
TGGTCAAAAA	TGATGAAATT	GACAAGTTTG	TCGCCCTTTA	CGAAGAAGAA	ATTATCAAGG	17760
ATAATACTGG	GAAACCTATC	CGTGTCAAAA	CCCTAGGGCA	AAAGCTTTAT	GTGGACAGTG	17820
TCAAACAGCA	TGATGTGACC	TTTGGAATTG	GGCCAGCAGG	TACAGGGAAG	ACCTTCCTTG	17880
CAGTGACCTT	GGCAGTGACT	GCCCTTAAAC	GTGGGCAAGT	CAAGCGAATT	ATCCTAACTC	17940
GTCCAGCGGT	GGAACCGGGA	GAGAGCTCTG	GATTTCTTCC	GGGTGATCTT	AAGGAGAAGG	18000
TGGATCCTTA	CCTTCGTCTT	GTTTACGATG	CCTGTATACA	AATTCCTGGG	AAAGACCAAA	18060

284

CGACTCGTCT CATGGAGCGT GAAATTATCG AAATTGCGCC CCTTCCTAT ATGCGTGGCC	18120
GGACCTTGA TGATGCCTTT GTCAATCTCG ATGAGCGCA AACACGACC ATCATGCAGA	18180
TGAAGATGTT CTTCGCGCT TTAGGTTTTT ATTCTAAGAT GATGTCAAT GGAGATATTA	18240
CTCAGATTGA CTTGCCACGT AATGTCAAGT CCGGTTTGAT TGATGCTCAA GAGAACTCA	18300
AGACATCCA TCAGATTGAC TTGTTCATT TTTCAGCCAA GGAATGGTT CGCCATCTG	18360
TTGTGCTCA GATTATCCA GCTATGAAT ATTCTACTGA AGTTGCACAC GACTGATTTT	18420
GAGGAAGTTC GCCTGCAAAA GAATAGACTT GTTCGTAAC TGTAAAAAT GTTATACTAT	18480
TTTTATGGAA ACAGTATACG ACAAGCACA AAACTTAAC TCAAAAACT TCAAACTATT	18540
GATTGGTGTG AAAAAGGAAA CCTTCAACT CATCTAGAA CACCTGAAT CAGCCTATCA	18600
GATTCAGCAC CGAAAAGGTG GACGTCCACG TAGTCTGCCC ATGGAAGACC AGCTCATTAT	18660
GACCCCTCGT TACTTGGCAT ATATCCAC TCAGCTCTG CTGGCCTTG ATTTTGGCGT	18720
CGGTGTAGCT ACGGTAAATG CCATCATCAC TTGGGTGGAG GATCACTTC GTGCCTCAGG	18780
TAGCTTTGAT TTGGACCAT TAGAAGCCCC GAGTCTGCT GTGGCTATTG ACGTACCTGA	18840
AAGTCCGATT CAGCTCCAA ACAAAACCAA AGCAAAAT ATTCTGGTAA AAGAAACGA	18900
CACACCTTAA AAACCTCAAT TATGCTGGAT TTGACGACAC ATAAAGTCTG TCAAAAGGCC	18960
TTTTCTGACG GACATACGCA TGATTTTACT CTCTTCAAAG AAAGTATTGG ACAAAATTTG	19020
CCTGAAACGA CGCTTGCCTT TGTGACCTA GCTATTTAG GCATCTGAA ATTTATGAG	19080
AATACTTTCA TTCTCTGTA AAATTCCAAA AATCGCGCGT TGAATGAGGA TGATAAGCAG	19140
TTAAATTAAG AGATGTCAGC GATACGAAT GAAATTGAAC ATTTTAACGC TAAATTCAG	19200
ACCTTCCAAA TCATGTCAGT CCTTATCGT AACCAGAGAA AACGTTTCA GTTACGGGCG	19260
GAATTAATTT GTGCCATCAT CAATTATGAA GTGAACCTAGA TTCCGAACAA GTCTAATATA	19320
CTTTGAGAG AGGAAAATCC AGTTGTATAG GCTAAAGCTT TTATCCAAAG GTCTGAGACA	19380
ACGATTAGGC ACGATGGAAA GAACTTTTAT GTGGCTGATG ACGATCAGTG CATCTTCTG	19440
TGTCATAATC ACAGGCGACA AGAAAGTAG AAATTGAAA GATGATTGAC CAATATCTA	19500
AGTATTACAG TTGTAGGATA CTAATGAAA AGGATATTC AAGTATTTTA TCTTTATATG	19560
AAAGTAATCC TCTGTATTT CAGCATTTGC CACGAGGCC AAATTTTCA ACTGTAAAAG	19620
AGGACATGCT TTGTCTAACC GAAGGTAAAG CTAAGGCTGA TAAGTTTTT GTTGATTTT	19680
GGAATGGATC TGACCTTGTG GCTGTTATGG ATTTTGCTA TGCATATCTT GATGAGGAGA	19740
CTGTTTTTAT TGGTTTCTTT ATGGTTGATC AAGCCTATCA GAGAAAAGG ATTGGTAGTC	19800
ATATTGTGAC AGAAGCACTA GCTTATTTTG CTAAGAAGCT TCGAAAGGCA CGTTTGGCTT	19860

ATGTTAAGGG AATCCGCAA TCTCAGCAT TTTGGGAAAA GCAGGGCTTT AATCAATTG 19920  
 GATCGGAGGT TAAGCAAGAA CTCTATACGG TTGTATTTCG TGAACAGAGC CTAGAAGATT 19980  
 AGAAGTGCA TCAAGTAAGA ACTATTGGA ATTGTGTTTG GAACAATTAT CAGGATTAGA 20040  
 TGATCTGACT TACCGTTCCA TGATGGGGGA GTATATTCCT TACTTCCGCG GCAAGATTAT 20100  
 TGGCGGCATT TATGACGATC GCTTTTATGT TAAACCCGTG CAAGCAGTCT TAGATAAGAT 20160  
 TGACCAATCT TCTTTTGAGT TTCCATACAA AGGTGCCAAA GAATTGATTT GAGTGGAAAG 20220  
 ACTTGATAAT AAGATGTTTC TATAAGACCT AATTTTAGCT ATGTATAACC AACTGCCAAC 20280  
 GCCCAACCT AAAAAGAAAA AGCAAGGGTG AACGAAGTAA AAGAAGATC TGCTAAGGCC 20340  
 CTCTCTTTGC ACGGGTAAAA TTTTATATAT AAAAAGAAGC TGGGACTAAA GAGCTCAGCT 20400  
 TCCTTTGGTT TATATAATTG TCATTACAAG ACGAAGTGGT TGGCGAAAC TCTGTTGACT 20460  
 TTATTTCAATT TAGAGTTCTT TATGCACAAT TGAGCTGGA ACGAAGTCT CCAATTGCAA 20520  
 AGTATACAGT ACAATAAAC AACGATGTAA TAGCTGATGA CACAAAGCAC AGTGGGTAGG 20580  
 ACTTGGCAAG TCACCCCTTT CTTTCAAJA TTTATACTAA ATCAATTGATA TCAGTGTAGT 20640  
 CACGATTAAG TCCTTGAGCA ACTGGTAGGT TAGTCAAGTA ACCTTGATAA GTAGTCACAC 20700  
 CTTGACGCAA GCCTTCACT TCAGAGATTG CTGTGCGAA TCCTTTGCAA GCCAAAGCTT 20760  
 CGATATAAGG AAGAGTGACA TTGGTTAGGG CGATGTTTGA AGTGCAGACA ACGGCACAG 20820  
 GGATATTGGC AACGCCATAG TGGAGAACAC CGTGTTTTC ATAGACGGGT TCATCGTGG 20880  
 TTGTTCACAG GTGAGCTGTT TCGATAACGC CACCTTGGTC AACAGCAAG TCAACGATAC 20940  
 AGAGCCTGGA CGCATTTGTT TGACCATCTC ATCTGTACC AATTCCGGTG CTTTTCAC 21000  
 AGGGATGAGA ATGGCTCCAA TCACCACATC AGCATCTCTC CACTTGCTT CAATGTTGAA 21060  
 TGAATTAGAC ATAAGGATT GAATTTGACT TCCAAAGACT TCTTCTAGAA CTGAGAGAG 21120  
 CTTGGAACATA ATATCTAAAA TAGTCATTG AGCACCAAGA CCAAGGGCGA TCGGGGCA7C 21180  
 ATGTCTACCG ACGACACCAC CACCGATGAT AGTTACTTTT CTTTTTGAA CACTTGGTAC 21240  
 ACCACCAAGT AGAACACCAG AGCCACCAGC TTGCTTAGTA AGGAAGTGAG CTCGGATTG 21300  
 AACAGCCATA CGACCTGCAA CTTCACTCAT AGGAACGAGG AGCGGTAGTT GTCTTGATT 21360  
 GTCACGAACA GTTTCAGTTG TTTTGTCTGT TAACATAGCA TCTGCTAATT CTGGAGCAGC 21420  
 GGCCATGTGC AAGTAGTGTA AGAGAAGAAG ATCCGTCGCG AAGTAACCGT ATTCAAGACT 21480  
 TAAAGATTCT TTTACTTTCA CAACCAACTC TGCTGCCCAA GCTTCACCAG CAGTAGCGAC 21540  
 AATCTCAGCT CCTTGCTTTT GATAGTCAGC ATCAGTAAAG CCAGAACCGA GACCAGCAAT 21600

	286	
TGTTTCGATA AGGACACGA? GACCACGACT AACTAAGCTA TGAACACCTG CAGGTGTGAG	21660	
GGCGACACGG TTTTCGTAT TTTTAATTTC TTTTGGGATT CGGATTACAA TTGAGATAAC	21720	
CTACCTTCA ATTGACGGTC TTGTTTGGT TGTACATTTC CAGTTCATAA ATCAAAAATG	21780	
TGACGGTTTC ATTGTATATG AAACCGCTTC AAAAAACAAG AAAAATTTGT CATCCAAATT	21840	
TTTTTATGCT AGACTAGTGA AAATCAAGCT CTAATGAGAG GAAAAGTATG GAATCAATAT	21900	
TTCTGAAAAT TGCCCATGAT CCGTCTATAG AAACGGAGCG TTTATTTGCTC AGACCTGTAA	21960	
CTTTGGATGA TGGGGAACAA TGTTTGACTA TGCCTCGAC AAGGGTAATA CACGTTACAC	22020	
TTTTCCRACC AATCAAAGCT TGGAGAAAC CAGAATAAAC ATTGCTCAGT TCTACTTGGC	22080	
TAATCCCTTG GGACGTTGGG GAATAGAACT AAAAAACAAT GGTGAGTTTA TTGGAACCAT	22140	
TGACTTGCAC AAGATTGATT CTGTTCTTAA GAAGGCAGCT ATTGGCTACA TTATCAATAA	22200	
AAAGTATGG AATCAAGGAT TAACGACAGA AGCCAATCGT GCTGTGATTG AGCTAGCTTT	22260	
TGAGAGATA GGGATGAATA AGTTGACTGC CCTTCAGAT AAGGCTAAT CCGGCTCAGG	22320	
AAAGGTCATG GAGAAATCAG GCATGCGTTT TTCCAATGCA GAACCATATG CTGTATGTA	22380	
CCAGCATGAA AAAGGCCGAA TCGTGACAAG AGTTCAATTAT GTCTTGACCA AGGAAGACTA	22440	
TTTTTGCAAT AATAAGCAG TTGAAAAGAA ATTTTTCGAC TGTTTTTCCT TCCTCTTACG	22500	
AATAATCTAA GAGAGGAGAA AATATGGAAG CAATTATCGA GAAAATCAAA GAGTATAAAA	22560	
TCATCGTCAT CTGTACTGCT CTGGGCTTGC TTGTAGGAGG ATTTTTCCTG CTAAAAACAG	22620	
CTCCACAAC ACCTGTCAAA GAGACGAATT TGCAGGCTGA AGTTGCAGCT GTTTCCAAGG	22680	
ACTCATCGAC CGAAAAAGAA GTGAGAGAG AAGAAAAGGA AGAACCCCTT GAACAAGATG	22740	
TAATCACAGT AGATGTCAAA GGTGCTGTCA AATCGCAAG GATTATGAC TTGCGCTGTAG	22800	
GTAGTCGAGT CAATGATGCT GTTCAGAAGG CTGGTGGCTT GACAGAGCAA GCAGACAGCA	22860	
AGTCGCTCAA TCTAGCTCAG AAAGTTAGTG ATGAGGCTCT GGTTTACGTT CCTACTAAGG	22920	
GAGAGAAGC AGTTAGTCAA CAGACTGGTT CGGGACAGC TTCTTCAACA AGCAAGGAAA	22980	
AGAAGGTCAA TCTCAACAAG GCCAGTCTGG AAGAACTCAA GCAGGTCAAG GGACTGGGAG	23040	
GAAAACGAGC TCAGGACATT ATTGACCATC GTGAGGCAAA TGGCAAGTTC AAGTCAGTAG	23100	
ACGAGCTCAA GAAGGCTCTCT GGCATTGGTG GCAAAAACA? AGAAAAGCTT AAAGACTATG	23160	
TTACATCGGA TTAAGAAATT CTCTATTCCC CTAATTTACC TGAGTTTCTT ATTACTTTGG	23220	
CTTTATACG CTATTTCTC AGCATCTTAT CTGCTTTTG TGGGCTTTGT TTTTCTGCTA	23280	
GTCTGTCTCT TTAATCCAAAT TCCGTGGAAA TCTGCTGGTA AAGTCTTAAT AATTGCGGA	23340	
ATCTTTGGAT TTTGGTTTGT TTTTCAAAAT TGGCAACAGA GTCAAGCGAG TCAAAATCTG	23400	

CGCGATTCTC	TTGAAAGGT	ACGGATTTTG	CCTGATACTA	TTAAGGTTAA	TGGTGATAGT	23460
CTATCTCTTC	GTGGCAAGTC	TAAAGGTCGT	GCTTTCCAAG	TCTATTATAA	ACTCCAGTCC	23520
GAGGAGGAGA	AAGAAGCCTT	TCAAGCTTTA	ACTGACCTGC	ATGAGATAGG	ACTAGAAAGG	23580
AAGCTTTTCG	AGCCAGAAGG	GCAGAGAAAT	TTTGGTGGCT	TTAATFACCA	AGCCTATCTG	23640
AAGACTCAGG	GAATTTACCA	GACTCTCAAT	ATCAAAACAA	TCCAGTCACT	TCAAAGATT	23700
GGCAGTTGGG	ATATAGGAGA	AJAAGTGTCC	AGTTTACGTC	GAAAGGCTGT	GGTTTGGATT	23760
AAGACGCACT	TTCCAGACCC	TATGGGCAAT	TACATGACAG	GACTCTTGCT	GGGACATCTG	23820
GACACGACT	TTGAGGAGAT	GAATGAGCTT	TATTCAGTCT	TAGGAATTAT	CCACCTCTTT	23880
GCCCTATCTG	GCATGCAGGT	AGGTTTTTTC	ATGAATGGAT	TTAAGAAACT	TCTCTTGCGA	23940
TTGGGCTTGA	CCCAAGAAAA	GTGAAATGG	CTGACTTATC	CCTTTTCCCT	TATCTATGCG	24000
GGACTTAAGT	GATTTTCAGC	ATCGGTATTT	CGCAGTCTCT	TGCAAAAGCT	ACTGGCTCAA	24060
CATGGGGTTA	AGGGCTTGGA	TAAATTTGCC	TTGACGGTGC	TGTGCTCTTT	TAPTCTCATG	24120
CCAAACTTTT	TCTTGACAGC	AGGAGGAGTC	TTGTCTTGCG	CTTATGCTTT	TATCTTGACC	24180
ATGACCAGCA	AAGAAGGGGA	GGGGCTCAAG	GCTGTACTA	GTGAAAGTCT	AGTCATCTCC	24240
TTGGGCAAT	TGCCCATTCT	ATCCTTCTAT	TTTGGCGAAT	TTCAACCTTG	GTCTATCTCT	24300
TTGACCTTTG	TCTTTTCTCT	TCTTTTGGAC	TTGCTCTCTC	TACCGCTCTT	GTCTATCTTA	24360
TTTGTCTTTT	CCTTTCTCTA	TCCAGTCAIT	CAGCTGAAT	TTATCTTTGA	ATGGTTAGAG	24420
GGCATTTATC	GCTTGGTCTC	GCAGGTGGCA	AGGAGACCAC	TTGTCTTTGG	TCAACCCAAC	24480
GCATGGCTTT	TAATCTTATT	GTTAATTTCC	TTGGCTTTGG	TCTATGATTT	GAGGAAAAAC	24540
ATTAAAGGAT	TAAACAGTAT	GAGTTTATTG	ATTACAGTCT	TCTTTTCTCT	TACCAAGTAT	24600
CCACTGGAAA	ATGAAATCAC	CATGCTGGAT	GTGGGGCAAG	GAGAAAGTAT	TTTCTACGGG	24660
ATGTAAGTGG	GAJAACCAAT	CTCATAGATG	TAGGTGGTAA	GGCAGAACTC	TATAAGAAAA	24720
TCAAAAAAATG	GCAAGAAAAAG	ATGACGACCA	GCAATGCCCA	GCGAACTTGT	ATTCCTATCT	24780
TCAAAGTCG	AGGAGTAGCT	AAGATTGACC	AGCTAATTTT	GACTAACACG	GACAAGGAGC	24840
ATGTTGGAGA	TTTGTACAGAG	ATGACCAAGG	CTTTCCATGT	AGGGGAGATT	CTAGTATCAA	24900
AAGACAGTCT	GAACAGAGAG	GAATTTGTGG	CAGAACTACA	GGCAGTCAA	ACAAAGGTGC	24960
GTAGTATGAT	AGTAAAGGAG	AAGTTGCCCA	TTTTTGGAG	TCAGTTAGAA	GTCTATCTCT	25020
CAAGAAAAAT	GGGAGTAGGA	GGACACGATG	ATACCCTAGT	TCTGTATGGG	AJATTCTTTG	25080
ATAAGCAATT	TCTCTTCAAG	GGAATTTTGG	AGGAGAAAGG	AGAGAAGGAC	TTGCTGAAGC	25140

		288	
ACTATCCAGA	CTTGAAAGTA	AATGTTTTGA	AAGCTAGCCA ACATGGCAAT AAAAAATCAT 25200
CAAGTCGAGC	CTTTCTAGAA	AAACTCAAA	CAGAGCTTAC TCTTATCTCA GTTGGAAGAA 25260
GCAATCGAAT	GAAACTCCCC	CATCAGGAAA	CATTGACACG ACTGGAAGGT ATCAATAGCA 25320
AAGTTTATCG	AAC TGACCAG	CAAGGAGCTA	TACGTTTAA GGGGTTGGAT AGTTGGAAAA 25380
TCGAAAGTGT	TCGATAGGAA	GGATAAATGT	TGTAGATTAG TGAATAAATC TAAAAATTTG 25440
TTGCATATAA	ATGATAAAAA	TGGTATAATG	AAAACGTATT CAATATTGAG GATATAAAT 25500
CATTAAAAAT	CAGCAAAAGT	TGTTTTATTA	GTTAGTTTAT AATCTATTGG TCTTCTTCAG 25560
TCAGGTGTAT	CTGCTGTGAC	AGTCACTAAA	AGTTACAAGT ATGATTGGAA TACGGTTTGG 25620
GAATATAGTA	CCAACTATCA	CGACCATCAG	TATGCTTGGA TTCCGTCATG GTCCTGTTAT 25680
GACAGCTATT	CTGAGTATAA	AGTTGGCGGA	GGCTGGAAC ACCTCGTTA TGAGGTCATA 25740
AAC TATTACA	GCGGAGGCTA	TTAATTTCTA	AAGAGTGAGA AAAAGGAACG CTAGATATGT 25800
TGCAGCTTAC	TCATGTGACC	TTAAAAACGC	GACAACTCAT CTTGCAAGT GTGGATTCCA 25860
CCTTTAAAAA	GCGTAGGGTT	TATGGTCTTC	TTGCTATCAA TGCTCTGGA AAGACGACCC 25920
TGTTCCGTGC	CATTAGCAAT	TTAATTCCCA	TAAGTAGTGG AAATATCGCA GCCCTTCCTT 25980
CTTTATTTTA	TTATGAGAGT	ATTGAATGGC	TGGATGGAAA CTTAAGTGGG ATGACTATCC 26040
TTGCTCTTAT	CAAAAACATC	TGGAGTCAG	GTCTGAACCT GAGGGATGAA ATCGCCTATT 26100
GGGAAATGTC	TGACTATATC	AGTCTTCCCA	TTCCGCAAGTA TTCTTAGGCG ATGAAGCAAC 26160
GCTTGGTGAT	TGCCATGTAT	TTCTTCAGTC	AGGCCAAATG CTGGCTCATG GATGAGATTA 26220
CAATGSCCT	AGATGAGTAT	TATCGACAGA	AGTTTTTGA TAGGCTAGCA CAATCGATA 26280
GACAAGAACCA	GCTGGTTCTT	TTAAGTTCCC	ACTATAAGGA AGAGTTGGTT GATGTCCTGG 26340
ATAGAGTAGT	AACCATTCAT	CAGGGGCAGA	TAGAAGAGGT TTAGTTTATG AAAGATGTTA 26400
GTCTATTTTT	ATTGAAAAAA	GTTTTCAAAA	CGCGCTTAAA CTGGATTGTC TTAGCTTTAT 26460
TTGTATCTGT	ACTCGGTGTT	ACCTTTTATT	TAAATAGTCA GACTGCAAAAC TCACACAGCT 26520
TGGAGAGCAG	GTTGGAAGAT	CGCATTCGAG	CCAACGAGAG GGCTATCAAT GAAAATGAAG 26580
AGAAACTCTC	CCAAATGTCT	GATACAGAGT	CGGAGGAATA CCAGTTTGCT AAAAAATAAT 26640
TGACGTGCA	AAAAAATCTT	TTGACGCGAA	AGACAGAAAT TCTGACTTTA TTAAGAAGAG 26700
GCGCGTGGAA	AGAAAGCTTAC	TATTTGAGCT	GCAAGATGA AGAGAAGAT ATGGAATTTG 26760
TATCAAATGA	CCCGACTGCT	AGCCCTGGCT	TAAAAATGGG GGTTGACCGC GAACGGAAAG 26820
TTTACCAAGC	CTGTATCCCC	TTGAACATAA	AAGCACATAC TTTGGAGTTT CCGACCCACG 26880
GGATTGATCA	GATTGCTCGG	ATTTTAGAGG	TTATCATCCC AAGTTTGTTT GTGGTTGCTA 26940

289

TTATTTTAT GCTAACACAA CTATTGCGAG AAAGATATCA AAATCATCTG GACACAGCTC 27000  
 ACTTATATCC TGTTTCAAAA GTGACATTGG CANTATCCTC TCTTGGAGTT GGAGTGGGAT 27050  
 ATGTAACCTG GCTGTTTATC GGAATCTGTG GCTTTTCTTT TCTAGTGGGA AGTCTGATAA 27120  
 GTGCTTTTGG ACAGTTAGAT TATCCCTACC CAATTTATAG CTTAGTGAAT CAAGAAATAA 27180  
 CTATGGGAA AATACAGAT GTATTATTTC CTGGCTTGCT CTTAGCTTTC TTAGCCTTTA 27240  
 TCGTCATTGT GGAAGTTGTG TACTTGATTG CTTACTTTTT CAAGCAAAA ATGCTGTCTC 27300  
 TCTTCTTTC ACTCATGGG ATTGTTGGCT TATTGTTTGG TATCCAAACC ATTCAGCCTC 27360  
 TTCAAAGGAT TGCACTCTG ATTCCCTTTA CTTACTTGCG TTCACTGGAG ATTTTATCTG 27420  
 GAAGATTACC TAAGCAGATT GATAATGTG ATCTAAATG GAGCATGGGA ATGCTCTTAC 27480  
 TTCTTGTCTT GATTATCTTT TTGCTATTGG GAATTTCTAT TATTGAAAGA TGGGGAAGTT 27540  
 CACAGAAAA AGAATTTTTT AATAGATTCT AGCTTTCTTA TAGGTAGGGA AATAAGTAA 27600  
 AAACCTAACAT AGAGAGGGAA TCAACTTGAT TCTCTCTTTT TGATTGCGAA ACCAAACCA 27660  
 AATACAACA CAACCTTTTC AAAAATAAC TTTTATCTTT GACAAAGCT AGAAAACCTG 27720  
 GTATCATATA AAAGTTGAGA AAAGCAGAAG TGAGAGCTTC TCGCTTTGTG ACATTAAGTT 27780  
 GCCTGGCCCT ACGGATGAAA AGTTTCGAAG AAACGCTATC ATAACTGCG GCTTTGTATA 27840  
 TTTACAAGTC CGCTATTGTT TTTCTCTAAT AAAACAAAAG AGGTGAAAAA CATAGCAAAG 27900  
 CAAGACTTAT TCATCAATGA TGAGATTCTT GTAAGTGAAG TTGCTTTGAT TGGCTTGAA 27960  
 GGAGAACAGC TAGGTATCAA GCCACTCACT GAAGCGCAAG CTTTGGCTGA TAACGCTAAT 28020  
 GTTGACCTAG TATTGATTCA ACCCCAAGCC AAACCGCCTG TTGCAAAAAT TATGGACTAC 28080  
 GGTAAAGTTCA AATTTGAGTA CCAGAAAGA CAAAAGAAC AAGCTAAAA ACAAGCCTT 28140  
 GTTACTGTGA AAGAAGTTCT TCTAAGTCCG G 28171

(2) INFORMATION FOR SEQ ID NO: 23:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 7147 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 23:

CCCTCAACT TTTGCAATCA AGGCTAAGTA GACAGCAGCA AATTCATAT TGTATAATTT 60  
 CTGACTATA CTTCTCTCTT TCTATGTGA CTAGTATAAA TAAGAAAAG AAGGCCCTCA 120

290

AGCCTTCCTTT	TGATTTATTTC	TTCTGCTTCA	TCTTCTGTAA	ATTGACTATT	GTACAAGTCA	180
GCCTAGAAGC	CACCTTCGCG	CATCAGTTCC	TCATAGTTGC	CTTGCTCGAT	GATATTTCCTA	240
TCTTTCATGA	CCAAGATCAA	GTCTGCAATT	CGGATGGTTG	ACAAGCCGTG	GGCAATGACA	300
AAGGATGTGC	GTCTTCCAT	CAAAAGGTC	ATGGCTTTTT	GGATCAATTC	CTCTGTCCGT	360
GTGTCAACAG	AAGAAGTCGC	CTCATCCAAA	ATCAAAAGCG	GTGCATCCCT	AAGAAGGGCA	420
CGAGCAATAG	TCAATAGTTG	TTTTTGTCTT	ACAGACAAGG	TCACGGTGTG	ATCCAAGATG	480
GTATCATAGC	CATCTGGCAA	GGTCATAATA	AAGTGGTGAA	TTCCCAAGC	CTTACTAGCT	540
TCCATCATTC	GTTCATCACT	AATCCCTATT	TGATTATAGA	TGAGATTGTG	TGGAATAGTT	600
CCTTCAAGAG	GCCAGGTATC	CTGCAAGACC	ATTGAAAGCG	CATCATGCAC	TTCTGAACGC	660
GTATAGCCT	TGGTATCCAC	ACCATCAATG	CGAATACTTC	CCTTATCAAT	CTCATAGAAT	720
TTTATCAAAA	GATTGACAA	GGTTGTCTTA	CCAGCCCCAG	TGCGCCCAAC	AATGGCAACC	780
TTTTGACCAG	CATGAGCTGT	CGCAGAGAAG	TCATAGTCTT	GAACATTGAC	ACCGTCCACC	840
AGAATTTTCT	CTGCTGACAC	GTGCTAGAAA	CGTGAATCA	GATTGACCAG	AGTTGATTTA	900
CCAGAACCTG	TTGACCCAAT	AAAGGCCAAT	GTTTGACCAG	TTTTGTCTTT	AAAGCTAACA	960
TGTTCAATAA	CTGCCCTCGA	ATTTGCCGCA	TAGCGAAGG	TCACATCCTT	AACTCGACC	1020
TGACCTTTGA	AGTTTTCATC	AGTCAGCTGC	ACTTGAACAG	GTTTGTGGAT	AGAAGAATGC	1080
AAATCTAAAA	CTTGATTAA	CCGCTTAGCA	GAGACCATAG	TTCCGGGAAG	AACGATGAAG	1140
AGTGCTCCCA	TGAGAAGGAA	GCCCATGACA	ACCTACATGG	CATAGACAT	GAAAACAATC	1200
ATGTCACATA	AGAGAGGCAG	ACGCGCTATC	GGAGCAGCGT	CGTTAATCAC	ATAGGCCCCA	1260
ATCCAGTAAA	TGCGCACACT	CAAAACCACT	GAAATCCCCA	TCATGATAGG	ATTCAAAATA	1320
GCCATAAGAC	GGTTGACAAA	CAAAATCAAA	CGGGTCAATT	CATCATTTAC	TGCTGCAAAT	1380
TTTTTCAATTT	GATAATCTCT	TGCATTGTAG	GCAGCAACGA	CACGAATACC	TGTTAAACTC	1440
TCACAGGTGA	TACTGTTCAG	TTTATCTGTC	AGCCCCGTAA	TCAAGGACTG	TTTTGGAAAG	1500
GCTAGCGTCA	TCAAAACGGT	CGTCATCAGG	ACGTTGATAA	TCACTGCCAC	AAGTACGGCC	1560
CAGAGCCAGT	ATTCTGAATG	ACCTAAAATC	TTCCCAATAG	CCGACATAGC	CATPAATTGAA	1620
CCACGCGTTA	CCACTTGCAA	GCCCATAGTA	ATCAACATTT	GAACCTGAGT	AATGTCATTG	1680
GTACTAGCGG	TCAAGAGGCT	AGGAATTGAA	AATTTCCTAA	TCTCTGTCTG	CGAGTAATCC	1740
AAAACTCGGT	TAAAAATATC	ACTTCTCAGC	CTACTAGTAT	AAGAAGCCGC	CACCTCGGAT	1800
GCAAAAAATC	CAACTGCAAC	TACGGACAAG	AAGGCAAGAA	AGGACATTCC	CATCATCATG	1860
CTTGGCGACT	GCCACAATCT	ATCTAAATTA	GTTCCTGTAC	TACCTAGCAA	ATCCGTAAAT	1920



TTCCGAGATAT	AGGTGGGCAC	TTCCAACTCT	AGATAGACCG	AAAAGCAAGT	AAAGAGAATG	1960
GCTAGTAAAA	TCATCCCCCA	TTCTTTTCTA	CTAATTCTTT	TGGCTAAITTT	CTTTATPCTC	2040
TCCTCTCTATT	CCCTTGATAT	TTTGCCGTGA	GTTGACCGAG	AACTCTCTCA	AAAATCAGTA	2100
ATTTCATCTTC	ATCAATGTCT	TCATCAACT	GCTTGCTCTAT	GCCTTCMAAA	AAAGCCTTAA	2160
CCTGTGTGCAT	CTGAGAAGCT	GCTTTGTCCG	TCAGACGAAC	AAACTTAGCC	CGCTTATCAA	2220
CAGGACTGGC	CTCCAATTCC	ACCAAACCAT	TTTGCACTAT	ACGCTTAACC	AGATTACTAG	2280
CAACAGGCTT	GGTAATATTG	AGTTCTCTGT	CGATATCTTT	AATCAAGACC	AAGTCTTGGT	2340
TTTTCTCGCG	ATTATCCAAA	AAACGCACAA	CCTGACCTTG	CGGCCACCC	ATAAATPCAA	2400
TGCGGCAACG	TTTGCTTCC	TTTTGCACCA	TCAGGTGAAT	TTGATGACCA	AAACGCTTAA	2460
AGACTAACAT	CGGTTTATCC	ATAATCTCCC	CCTTCTAAAT	AAAAATAGTT	CTCTGGAGAA	2520
TAAATTAAAT	TCTATGAGAA	CTATTTTCTT	GATTAACAAA	ATCCCAAGTG	ATTTTCTCAC	2580
TTAGGATCAT	GTTCTATAGG	TAAATTAATA	ACCCATCTAC	GTTCTATATA	ATCTTTTGGA	2640
CGTCTTGTCT	GTTCTCAAGA	ACGCTGTAAA	GTTTTTCAAA	GGTTTCAAGG	TCTTCGCCGT	2700
ACAATTCCAC	TTCTGACTGA	GGAATCATTT	CCAATTCAGT	CACTTGGAAT	TCTTCAATAC	2760
CAGACTCACG	GAGGGCAACG	ATAGCCTTGT	GAAAGTCAAT	TGCGCTGTGT	TAAAGCTGTA	2820
TTGTACCTTC	TTGTGCTTCT	AGGTCAATCCA	CATCCACATC	CGCTTCGAGC	AAATTGCTCAA	2880
AGACTGCGTC	CGCATCTTCA	CCTCCAAATA	CAATAACACC	TTTTGTTGTCA	AAGAGGTAAG	2940
AAAACAGAAC	TGAAGCGCCC	ATGTTTCCGC	CGTTTTTACC	AAAGGCTGCA	CGGACATTGG	3000
CTGCTGTACG	GTTGACGTTA	GAAGTCAAAG	TATCCACAAT	TAGCATAGAG	CCATTTGGCC	3060
CAAAACCTTC	GTAACGTCTT	TCGTAAAGG	TTTGTCTCTG	GTTTCTCTTG	GCTTTATCAA	3120
TGCTTTTATC	GATAATGTGT	TTTGGCACTT	GGGCTTGTGT	AGCACGGTCG	ATAAGGAATT	3180
TCAAAGCTGA	GTTTGATCTT	GGATCTGGAT	CACCTTTTTT	AGCTGCTACA	TAGATTTCTA	3240
CACCAAAATT	TGCATATACT	TTAGAGTTAG	CTCCATCTTT	AGCCOTTTTC	TTGGCTACGA	3300
TATTGGCCCCA	TTTAACOTCCC	ATTAGGAATC	TCCTTTTTTC	ACAATTTAAT	CTTTCTTATT	3360
ATAACACAAG	TTTTTTTTCAT	TTTCACTAGA	GGAAATGGAT	TTTATTAGCA	AATCAAGCTA	3420
GGATAGCACT	TTACCTGCTA	AGATGGTCTT	GCCTTTCTAT	CTTTATCAAC	AGGCACCTCAT	3480
CCAACATTCAA	AAAACAAACT	AGACCAATTAT	CTGCATAATAG	AAAGTTTCAG	CCAAGTTTGA	3540
CAAAGTCAGC	TCAAATTAAT	GTTTGAAGTT	TGTAGATATA	AGCGACAAAA	ACAATCATAC	3600
TGCACCTTTT	GTTGACAGTC	TACTCCAGAC	ATATCATAGT	TCAAGTAAAT	ACTTTGAAT	3660

292		
TCACAGTTC TTAAGGCGC TAATGTATTC TAAGAAATCA ATAGAAGAGT TTCTAAGCAA	3720	
ACCTCTAATA CTCAATAAAA ATCAAAGAGC AAATAGAAA GCTAGCCTCA GGTGCTCTCA	3780	
AACACTGTTT TGAGGTGCG GATGGGGCTG ACATGGTTTG AAGAGATTTT CGAAGAGTAT	3840	
AATTTAAGTG TTCCAAGAT GGAGAAGTTA GACTAGTACA CTGGCACTTC TAAAACATTG	3900	
CTAGCAATFG ATTGTTCAT ATTTAATTC ATTTTTTCCA TAAATGGGTA TTAGATATAA	3960	
ACAGCAAAAT ATTTCGATA CGTGTGCTC TTGAATTTCC AATCATCTAA AACAGTAAA	4020	
CGATAATCAA TCCCCTGTAT ATCAAGGAAT TGGCTACCT TTTTACTTTT TTACACATTC	4080	
TGTTTGATAG ATTCAATTTA ACATCAGGAG CATACTCCAA TGGAAATCC TAGGCAAGAG	4140	
ATAAATTTTC AGATATCCGC AGAGAGATCA TCGCTCTTT TTGTCGCAAG CATTCCTCTC	4200	
TCCTAGTCAT TTTCTACCTT ATCTTCTACC TGAGGATAGA GAGTTGTTCC CCAATAGAA	4260	
ATCGTCCGCT TACGCACTAG TGGCAAATCG GTTTTTTCAT AAACCGTAGC CCACCATTC	4320	
CAGGCAAGCC CGGTACACTC TCTAATTTTG ACAGAGAGAT TACGAACATP CCCTTTTAAA	4380	
GGAACTACTAG TGCTAAAGTG AGCGGTTAAA TCCTGCGCAT TTCTGTCCA AGCCTTAGGA	4440	
GTCAAGACTT CCTTACCTTG ATGATCATAG GATAATTCAT TCCAAATAAT ATAATATTGG	4500	
GCAACATAGG CACCACTATG ATCCAGCAGT AAATCTCCGT TTCTGTAAAG TGTAACTTTA	4560	
GTCTCAACAT AGTCTGACT ATTTTGAAG GTGCAACTA CATGTGACG TAAAAAGAA	4620	
GTGTATAGG AAATCGCAA GCCTGGATGA TCTGCTGTAA AGGCACTGCC TTCTTGAATC	4680	
AAGTCTCTA CCATATCCAC CTGCGCTGT ACACTCGG CACCCGAAT TGGGTGCGCC	4740	
CTTAAATAA CGGCTTCCAC TTCTGTATG TCCAAATCT GTTTCACATC TGCTGAGGA	4800	
GCTACCTTGA CTCTTTTAT CAAAGCTTCA AAAGCAGCCT CTACTTCATC ACTCTTACTC	4860	
GTGGTTTCCA ACTTGAGATA GACTTGGGCG CCATAAGCAA CACTCGAAAT ATAGACCAA	4920	
GGAGCTCTG CAGAAATTC TCTCTGTTTT AAATCCTCTA CGGTACAGT ATCTTGAAAC	4980	
ACATCTCCTG GATTTTTAAC AGCATCTACG CTGACTGTAT AATAAATCTG CTTAAAAATTA	5040	
ACAAATCGAA TCTGTTTTTT GCCTGAATGG ACAGAGTTAA AATCAATATC AAGAGAAATC	5100	
CCTGTCTTTT CAAAGTCAGA ACCAAACTTG ACCTTGAGTT GTTCCATGCT GTGAGCGGTG	5160	
ATTTTTTCAT ACTGCATCT AGCTGGGACA TTATGACCT GACCATAATC TTGATGCCAC	5220	
TTAGCCAAACA AATGCTTTAC CGCTCGGCA ACACCTGAAT TGCTGGGCTC TTCACTTTGG	5280	
AGAAAGCTAT CGCTACTTGC CAAACCAGGC AAATCAATAC TATAAGTCAT CGGAGCAAGA	5340	
TCGACCGCAA GAAGAGTGGG ATTATTTCTCT AACAGGTCT CTCCACTAC GAGAAAGTCT	5400	
CCAGGATAGA GCGCACTGTC GTTGGTAGCT GTTACAGAAA TATCACTTGT ATTTGTCGAC	5460	

AAGCTCGCT	TCTTCTTTC	GATAACAACA	AATCATCGG	GTAGCTGATT	ACCCCTCTTTG	5520
ATGAACAGAT	TTTCAATACT	TCTCCCTGA	TGGGTCAAGA	GTTTCTTTT	ATCGTAATTTC	5580
ATAGCTAGTA	TAAAGTCATT	TACTGCITTA	TTTGCCATCT	TCTACCTCCT	AATAAGTTCC	5640
TGGATTGAGT	TGCATAAAT	CAGACTTGTT	CAGCGAAATC	AGCCGTGTT	GGACTAAGTA	5700
ATCCAAAT	TCCCTGTACA	ATTCTTCTGA	GACATTCGT	CGCCGTCTG	CTAAATAAGA	5760
AGTGGAAATG	ACCGTATTAT	CCAACATAAA	TACCTTATCT	AAGTCAATCA	AGGTTGGTCT	5820
TGTAAAAGGA	TTACGAGCTA	GATCCGGCTC	TTCTATCATA	AAGTTCTTGA	CCAAAAGTCT	5880
GGTCAAGAGA	GCTGGTTTGA	AGGTCTGATT	TTTAACCAAC	TCTTTGTTTT	TAGTCATGCT	5940
GTTGTCAATA	CAGATATACA	TATGATTCTT	CACAGCCAAA	TGCTACTTAA	TAGTCGGAAA	6000
AGGCAATATA	AGAGCTACAA	CATCTCCTCT	CTTAATCAAG	CAAGAGCACC	CCCTTTTCTC	6060
CTAATGTAA	ATAGACAGGA	TTGACCAAGT	CTTCTGATTG	ACTCAGAAAT	TCCAAAGTTT	6120
GAGTTTGGCG	CGCTGTCAAT	TTAGTAGCAT	CTTGCTCTCT	CAATACAAAA	TGCTTGTCGC	6180
CAATAACCTT	GACAAATATA	TCTTCTCCA	AAGCTGACTG	GTAAATCCAC	ATCAGATGTT	6240
GTCTGTCTCG	AGAACTCAAG	AGAGAAGGAT	TTTCAAGCCT	CCCGATAGTC	TGATAAAAAAT	6300
CAAAAACAGG	AGCTAACTCC	TGCCAATCTG	ATTGGCTAGT	TGTCAGGGCT	AGAAAAAGG	6360
CTTTGCGAGC	TGATACTTCT	TGGTTAGCCT	TGAGAGTTAC	TTTCCCTCC	AGTTTTTATA	6420
GAAATCGGGA	AACTCCAGAA	AGCAAAATTT	TCTCTAACTG	CGAGAAATAA	AAACCTTTTCG	6480
TTCCCAGACA	TAAGTCTTTC	ATGTCGCTTT	CTCTAGCAAA	TAAGAGCTCA	AACATTTGAT	6540
AGTAAAGAA	AAATATCTGG	CACCTGGTCTG	CGCTCATCTT	TTCTTATCG	GCTTCTTTTT	6600
TTAACCCAGAG	CAAGGGCGAC	AGGTAGCTGG	ATTGAGACAT	TTCTCTTACC	TCTTACTCTT	6660
TTTTAACTGG	AGCATCTGCA	CTAGCTGCCA	CTTCTTTTGA	CTGGATACTT	TCCCCTGCT	6720
TAATCTCTTC	TGAGATAAGA	CTTTCGCATG	TCCTGACAAA	TAGGGCAAAA	GCCTTGGCTT	6780
TTCTCTGCATA	TTTCTCCGTT	TGGCATTTGAT	AGAGGAATTT	TTCTTCTCT	AGGAGTTGCG	6840
CAGTTTTTTG	GTAAGAAATC	CAATTTTCTT	TTGCATTATA	CAAAATGATA	ATCCCCCTCAC	6900
ACAGCAAGCC	GAGACTGGAT	AAGCAACCG	AAATCAACAG	GTAGCGATCA	CCTGGCATAG	6960
GAATAGACACA	AAAGACAGCT	ATGAGGAAAC	CTGCCACGAT	TTCTGTATT	TTTAATACCT	7020
TATAGCGCCT	ACGATGTTGA	ACGCTTTTCT	TTAAAAATG	AGCTATCTGT	ACGCTTAATC	7080
GCTCTGTCAG	GTACATTTCT	TCTGGCGTCA	TATTCGTAAC	TCCTTTTCTT	TACTTTTGATA	7140
ATCAGGG						7174

294

## (2) INFORMATION FOR SEQ ID NO: 24:

## (i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 755 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 24:

```

CCGCATGGGA TTGGTGCCT TTGGGCAAT CTCCTTGACC AAACGGAAA CATGTTTAT      60
GCGCCTGCCT TTACTGCCCT TGTGGGGT ACGTCTATAT GATCCTAGTC GCAAAAGTTC      120
CGCGCTTTGG AGCCATTACC ACTATCGGCC TTGTCATTGC CTCTTTTTC TTGGGAACATA      180
AACACGGTGC TGGTTCCTTC CTCCTGGAA TTATCTGTGG CCTCCTAGCA GATGGAGTAG      240
CTCATTTAGG AAAATACAAG GACAAAACAA AGAATTCCT TTCTTTCATT ATTTTCGCT      300
TTAGTACAAC AGGACCAATC TTCTTATGT GGATTCGCG CAAGCCTAT ATGGCTACTC      360
TTCTGGCAAG AGGAAAATCC CAAGAATATA TCGACGTAT CATGTCGCT CCAAAACCTG      420
GAACTGTCCT TCTATTATC GCAAGTATTG TCATCGGAGC CTAAGTGGGT GCCTTGATTG      480
GACAAGCCTT GAGTAAAAAA TTGCCCAGA AATCTGATC AGTTAAAAAG AGCCACGCGG      540
CTCTTTTATA TTTATGGCTC AATTCTTAG TCAAGAAAT TCCCAAGAAT TGGATTGCAA      600
AGATAATCAA AATGATAATA ATGATTGCCA AGATGTGAC ATCGTGATTG TAGCGGTAA      660
ATCCATAAGC GATGGCTAGC TTACCGATAC CACGAGTCC AACCGCACCG GCCATAGCTG      720
TTTCCCAACA AGGGAATCAA GGTACAGTC GTAC                                     755

```

## (2) INFORMATION FOR SEQ ID NO: 25:

## (i) SEQUENCE CHARACTERISTICS:

- (A) LENGTH: 3010 base pairs
- (B) TYPE: nucleic acid
- (C) STRANDEDNESS: double
- (D) TOPOLOGY: linear

## (xi) SEQUENCE DESCRIPTION: SEQ ID NO: 25:

```

TTCAATTGGT ATCTCAATCA ACGGTCTTCA CATGGTTTCA ACTGGTTTGA CTCTTGAAAA      60
AGCGAAAGCT GCTGGTTACA ACGCAACTGA AACAGGCTTT AACGATCTTC AAAAACACAGA      120
ATTCAAGAAA CATGACAACC ATGAAGTAGC AATTGAAGATT GTCTTTGACA AAGATAGCCG      180
TGAAATTCCT GGTGCCCAAA TGGTTTCACA TGATATTGCA ATTAGCATGG GAATCCACAT      240
GTCTCATTCT GCTATCCAAG AGCATGTGAC AATTGATAAA TTGGCATTGA CAGACCTCTT      300

```

CTTCTTGCCA CACTTCAACA AACCATACAA CTACATCACA ATGGCTGCCC TTACGGCTGA	360
AAATTAAAAA TGAATGAGCT ATCTGGCCTT AAGTTAAGGT CAGATAGTTT TTAGCTAATT	420
TGTCCCCATA CAATATAGT TTTTTFATCT TGTGCTTCAT TCTGTTCTGA CTTAAATGA	480
AAAGTGTAGT ACCAATACAA ATGATGAGGA TAAACAAT GACTGAAAT CGTTATGAC	540
TAAATAAAAA CTGGCACAG ATGCTCAAGG GTGGTGTAT TATGGATGG CAGAATCCTG	600
AACAGGCTCG TATCGCAGAA GCTGCTGGTG CGGCAGCTGT GATGGCCTTG GAACGAATTC	660
CGCGTGATAT TCGTGACGCT GGAGGAGTTT CCCGCATGAG CGACCAAAAG ATGATTAAAG	720
AAATCCAAGA AGCGGTTAGT ATTCCAGTAA TGCTTAAGGT CAGAATCGGG CATTTGPTTG	780
AAGCTCAGAT TTTAGAGGCT ATTGAAATTG ATTATATCGA CGAGAGTGAA GTTCTATCTC	840
CAGCTGATGA CCGTTTCCAT GTGGACAAGA AAGAATTCCA AGTTCCTTTT GTCTGTGGTG	900
CTAAGGATTT GGGTGAAGCC TTGCGTCOTA TCGCTGAAGG TGCTTCCATG ATTCTGACCA	960
AAGGAGAACG AGGACAGGG GATATCCTCC AAGCTGTTGG TCATATGCGT ATGATGAATC	1020
AGGAAATTCG CCGCATTCAA AACTTACGTG AGGACGAGCT TTATGTTGCT GCCAAGGATT	1080
TGCAAGTCCC TGTAAGATTG GTCCAAATATG TTCATGAACA TGGAAAAATG CCAAGTTGAA	1140
ATTTCCGTGC TGGAGGTGTT GCAACGCCAG CAGATGCTGC GTTAATGATG CAATTAGGGG	1200
CAGAGGGGGT CTTGTGCGT TCAGGTATTT TCAAGTCAGG AGATCCTGTT AAACGAGCGA	1260
GTGCCATTGT TAAGGCTGTG ACTAAGTTCC GTAACTCTCA AATCCTAGCT CAAATCTCTG	1320
AAGATTTAGG AGAAGCCATG GTTGGTATTA ATGAAATGA AATCCAAAT CTATGCGCTG	1380
AACGAGGAAA ATAGATGAAA ATCGGAATAT TGGCTTGCA AGGGGCCCTT GCAGAACATG	1440
CAAAAGTGCT AGATCAATTA GGTCTCGAGA GTGTAGAAT CAGAAATCTA GATGATTTTC	1500
AGCAAGATCA GAGTGACTTG TCGGCTTTGA TTTTGCCCTGG TGGTGAGCTC ACAACCATGG	1560
GCAAGCTCTT ACGTGACCAG AACATGCTAC TTCCCATCCG AGAAGCCATT CTATCTGGCT	1620
TACCACTGTT TGGGACCTGT CGGGGCTTAA TTTTGCTGGC TAAAGAAATC ACTTCTCAGA	1680
AAGAGAGTCA TCTAGGAAT ATGGATATGG TGGTCGAGCG TAATGCTTAT GGGCGCCAAT	1740
TAGGAAGTTT CTACACGAAA GCAGAATGTA AGGGAGTTGG CAAGATTCCA ATGACCTTTA	1800
TCCGTGGTCC GATTATCAGT AGTGTGTGGT AGGGGTAGA AATTTTAGCA ACAGTGAACA	1860
ATCAAAATTG TGCAGCCCAA GAAAAAATA TGTGCTAAG TTTCTTTTCA CCAGAAATTGA	1920
CTGATGATGT GCGCTTGAC CAOTACCTTA TCMATATGT TAAAGAAAA AGTTGAGATT	1980
GAATTTCTCA ACTTTTTC ATGTAATAAA CAATAGCGAT GTATTGAAGT GCGGACGACG	2040

296

CTAGGATGAA GAGATGCCAA ATCATGTGGA AATTAAGGTTT TTCTTTGGCA TAAATCCAG	2100
CTCCAACTGT ATAACAGAGT CCGCCAGTTA CCATGAGACT CCAGAAAACG GGTGTGCTTT	2160
GACTGATTAAT GGCAGGAATG ATAGCCAGAA CCAACCAGCC CATTAATCAGG TAAAGAGCAA	2220
GGCTAAATTT CTCATTGACC TTTTAGCAA AGATTTTATA GAGAATACCA AAGATGCTCG	2280
TTCCCATCTG GATGACAAAT ATCAGATAGC CAAACAGTTP ATTTCATCAAG GTCAAGACAA	2340
CGGCGCGTGA TGAGCCCGCA ATGGCAACGT AAATCATAGA ATGGTCAATG ATTGCGCAAAA	2400
CATATTGTG GGTGGAACCA TAGGCCATAG AGTGATAAAT GGTGGATGAT AGGAACATGA	2460
GAAAGAGACT GATGACGAAA ATGGAACGCG CGATAGAGGA TAAAAATCTG TGTGCTTCAT	2520
AACATATAGT GGATGAAATA GGCAGCAAGA TAAGCATGAT GACTGCACCC ACAGCATGGG	2580
TCACGCTATT AGCAATCTCC TCTCCAAAAC TGAGTTGTTT GCTGAGTTTA AGACTAGTGT	2640
TCATTTGGATT ACCTCTCTCT GAGTATGATC GATTAACTCT AGAGTTTGTG GATAGAGTTT	2700
AACGGTTTGG CAGCTGGTTT GATATAAGG GTTAGCTGGG TCAATTCTCT GGTTCATGTA	2760
GTCCACAAAA GCATCGTAGA GTTGGTCTGA ACTTGCTTGA GTTTGTAGAG TATTAAAGTGT	2820
CTGGGCTATT TCTTGAATAG AAAATACAGA CTTGAGGGTT GTGATAGCAA TCAAACGGGC	2880
AATCTGTTGG CGTTGGTATT TTTTTTTGTC AGGCTTTGTC AGGTAACCAT TTTTCACATA	2940
ATTGTTGACC ATAGATGCTG TTAGGCCCTT GTCTTTATTA GGAGAGATAG GGGCGCAGAC	3000
CTGATTGACA	3010

(2) INFORMATION FOR SEQ ID NO: 26:

- (i) SEQUENCE CHARACTERISTICS:
- (A) LENGTH: 15213 base pairs
  - (B) TYPE: nucleic acid
  - (C) STRANDEDNESS: double
  - (D) TOPOLOGY: linear

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 26:

CATAAATCGG TGCAAAATAC TTAATAGTGA AGTAGCAATT TCTTTGCTAT TTACTGAGG	60
CATATTCCCT AGACGAAAGA ATATTATTAT CAATCAAACT ATTGAATGAA CGTAGCTTTT	120
CAACTCTTC TACTGTTAGA TTTCTGACAA CATTTGTTGC ATAGACCTTA TTTCCATCAG	180
GATCAGGATG GTACTCATTT GTAACTTTTC TAAGAAGTTG TTGTTTTTGA TTTGATATCCA	240
ATTTAAGAAT TGAATTTCTC TCGAGATATT CCAACATATA AACCAACGTC AACATGTTGT	300
GGACATATTG CTTCAATCA TCTGCATTAT TAAATCTTGT AGTTGGATCA AGTACTTGTA	360
ATCGTCACT TTCTGTACTA TCAGATTTTG AATGTTTCAA GATGAGATTG ATGGTAATGG	420

TCGCATCATC TGGATGGTCT GGTGCTTGTA ATAACTCTTT AGCAAAGAAC TCTGGTCCCA	480
AGCCACTTCT TCGACCATAT CCTCCAAGAT AAATGTCTCG ATCTGAGTCA TGTGTCAATCT	540
CATGCGTATA AGTAATAGCT CCATCCTTAT CCAACATTGG ATAAACCCATA TAATAAACTG	600
CATCACTCTG AGCATAAGCA CCGTGTGTGAT TATGCCCAAC TTTATTTCOA ACAGGTCCAA	660
AGAAATGTGG CATTGCAGGA TTTGGATTAT CAAAATCTGC CACTTCTGTA GCTTTCCCTA	720
CGGTATTATC ATCGCCAAAT TTATAAGCAT CGTAAGCAA AATATTCTTA TAAAGTTTTT	780
CACGTGCATT GTGCTCTAAA ATACGATACC AATAATCGTA GTGATCTGCG TGACGTTTGG	840
CTGTTTCACG CGCATTTTCT TCAACAAAAT CATTGAGAGC CTTGCCCGCT TTATGGTCAC	900
TACTGGGATA CGGATCATAA GCTCCAAATC CTAGACTAGA CATGTCGAG ATGACAAAATA	960
CGGATCTCTC TGGCAAGGTC AGGAGAGGCA AGACCATATT GCGGTATTTC CATGTGOCAC	1020
TCGTGATPAG ATCATAAACA CCGATAGAAT ACTTGGTGGC AGCTAACCTT TGCTTCGTTT	1080
TCACCTCTTC GATAGTGAT TTTTCTTCGA CAATGTAAGC CTTAGTCTCT GATTTAAACC	1140
AGTCATTATT GCTTGTATTT GGTAAAAGA CTTTTCGGTA ATGTTCCAGC GTGCTAAACA	1200
AACTGTGCTG TCCATGTGTA CTGGCAAGAC TGATACCATA AGTATCGACA TTATTCTTAG	1260
CTAGAAGATT GTTAAAGCCA GATTTACCCA ACTCAATCAG AGTATCTAAT GGTGAAGCAT	1320
TCCCTTTACC AAAGAAGTCC AAATGTTACA GAACCTAGGC TTTGACATTC ACCTGACCAT	1380
AGCTAAAGTT ATACCAACCGT TCCAGATAGG TCAAGCCAAG TAGCAAGGCT TCCTTGTGTC	1440
GTTTGATTTT ATCTACAAGA TAACCTTCAG TGACGGGGTT AGCACTAGCC AGTCCAGCAT	1500
CCGCTGACAA GAGTTTTTTC AAACGTCTCT CCAGTTGTTG TTTTGTGTTG GCGAACTGGT	1560
CTTCTGATATA GAGCTCAGTT TGCCTGACGT TTGGAGAAAT ACCCAGCGTC TTTCTGATGG	1620
CTTCTGAATG ATAGTCAACC TTTTGTAAAT CAGGTAAAGC TTGCTTGATG ATAGAGGTTT	1680
GGTCATACAG GAATTGGTTT GGCCTATAGA GAAGTCCAGT ATTGCCCAGA CTATATCTTG	1740
CTAAATTTGGC GAAATCAATTC TGTATTGGA GATCCAGCTT CTCAGATAAA TCATCTCTGT	1800
AGTGAAGCAA GAGTTTGTGT GCAGTCTGTT TGTAGAAAAC AATGTCTGTG ATGACTTGTT	1860
TGTCTTCAAT CATGACTGCT GACAAGAGTT CTTTTTGATA TAAAGACTG TTTCTAATGA	1920
CCAGTTCTCC GTATTGACG ATGTTTGCTT TGTGTAGAA AGGTGACAA TTTTCAATGT	1980
TTTTATAAGT CAAGTTGCGC TTAGCTTGAT AATAGGCCAC CTTAGAAAAA TCACTGTCTT	2040
TTTTGCCACT TGTTGAAAGT GCTCCACTG TTGTTAAAA GAGAGGATTT ATTTCTGCTT	2100
TTTTGCTTGC AATTGAGAA GCATCTAGCA TTGTCCTCTT TTCTTCAAG GATTCCTTGC	2160